

الجينوميكس والمعلوماتية الحيوية

Genomics & Bioinformatics

رئيس التحرير
أ.د. أحمد يوسف المتينى
كلية الزراعة - جامعة الإسكندرية

المشاركون في التحرير
أ.د. سناء أحمد رياض
أ.د. آمال أحمد عبد العزيز
كلية الزراعة
كلية الزراعة
جامعة الإسكندرية
جامعة الإسكندرية
مهندسة الوراثية
مدينة السادات جامعة المنوفية

أ.د. ياسر محمد مبروك
أ.د. أحمد محمد الشاوى
كلية الزراعة - جامعة الإسكندرية
كلية الزراعة - جامعة الإسكندرية

م. أمير السعيد يسن
كلية الزراعة - جامعة الإسكندرية

2006

مكتبة بلستان المعرفة

طباعة ونشر وتوزيع الكتب
كفر الدوار - الحدائق بجوار نقابة التطبيقيين
١٢١١٥١٢٣٧. & ٠٤٥/٢٢٢٤٢٢٨ : ☎

•

•

•

•

•

•

الجينوميكس
والمعلوماتية الحيوية
**Genomics &
Bioinformatics**

الرجـاء مـن و المـعلـومـاتـيـة الـحيـويـة

أ.د. أحمد يوسف المتـرنـى

٢٠٠٥/ ٢١٩٠٢

I.S.B.N 977-393-057-2

الأولى

مكتبة بلستان المعرفة

كفر الدوار الحدائق ٦٧ ش الحدائق بجوار نقابة التطبيقيين

☎ : ٠٤٥/٢٢٢٤٢٢٨ الإسكندرية ٠١٢١١٥١٢٢٧

Email: bostan _ elma3rafa @ yahoo.com

اسم الكتاب

اسم المؤلف

رقم الإيداع

الترقيم الدولي

الطبعة

الناشر

جميع حقوق الطبع محفوظة

ولا يجوز طبع أو نشر أو تصوير أو إنتاج هذا المصنف أو أى جزء منه

بأية صورة من الصور بدون تصريح كتابى مسبق.

تمهيد Preface

يعتبر الكثير من المفكرين أن اكتشاف التركيب الجزيئي للمادة الوراثية (الـ DNA) من أعظم إنجازات القرن العشرين قاطبة، فبه وما تلاه من دراسات أصبحنا نعرف لغة الحياة. فالنيوكليوتيدات الأربع هي كل حروف تلك اللغة، تتابع في ثلاثيات كل منها يشفر لحمض أميني بعينه، ومن تلك الأحماض تتكون توليفات لانتهائية من البروتينات عصب حياة الكائن. فإذا ما اعترى المفكرين ذلك بأنه إنجازا إنسانيا، فهو لنا نحن الوراثيون ثورة قلبت علوم الحياة ودفعتها إلى أفاق لم نحلم بها من قبل ومع ظهور تقنيات الهندسة الوراثية أصبحنا قادرين على تطويع الوراثة تبعا لرغباتنا، وأصبح من الصعب على المتخصصين أن يلاحقوا أو يتابعوا مكتشفات علوم الوراثة الجزيئية خصوصا والبيولوجيا الجزيئية عموما، وأصبح البحث في تلك المجالات كل يعمل بعيدا عن الآخر مثل الجزر المتباعدة في المحيط، مما أفقدنا القدرة على استخلاص الكليات وغرقنا في بحر التفاصيل الصغيرة غير المترابطة، وظل الحال على ما هو عليه إلى أن قامت مجموعة من علماء الأحياء والرياضة في مطلع عام ١٩٧٩ في اجتماع جمعهم في أحد قاعات جامعة روكفلر Rockefeller University بمدينة نيويورك حيث دعوا إلى إنشاء قاعدة معلومات لحفظ تتابعات الـ DNA، ومن بعدها مباشرة وفي السنوات القليلة التالية ظهر اتجاه جديد لمحاولة الربط ما بين علوم الحياة وعلوم الكمبيوتر والبرمجة ويعرف هذا المجال باسم المعلوماتية الحياتية Bioinformatics، هذا المجال الجديد أدخل تغيرات جذرية في المناهج العلمية لدراسة علوم الحياة خصوصا البيولوجيا الجزيئية، وقد تسابقت الجامعات والمعاهد العلمية في دول العالم المتقدم إلى إضافة هذا المجال الجديد إلى مناهجها حتى تواكب التغير المأمول. ونحن هنا لسنا في صدد دراسة الـ Bioinformatics من حيث أسسها النظرية، فذلك يندرج تحت علوم الكمبيوتر والهندسة، بل يهمنا أن نعلم هنا التطبيقات والتقنيات التي يمكن أن تفيدنا في دراستنا للوراثة الحديثة. ويمكننا أن نسوق هنا، مثلا طريقا يؤكد أهمية المعلوماتية الحياتية في مجال العلوم البيولوجية الحديثة. فقد اختير في مطلع عام ٢٠٠٠ طالب الدراسات العليا James Kent الذي يدرس في معمل العالم البيولوجي الشهير Zahler بجامعة كاليفورنيا ليكون بطلا قوميا بواسطة جريدة New York Time لتجازه في عمل برنامج للكمبيوتر أطلق عليه اسم

”GigAssembler“. هذا البرنامج مكن الفريق القومى من مشروع الجينوم البشرى من حل تتابعات حوالى ٤٠٠,٠٠٠ قطعة متداخلة من الـ DNA (contigs) كانت مازالت باقية لم تحل بعد، هذا البرنامج أتمها فى غضون أربعة أسابيع فقط. هذا الإنجاز مكن الفريق القومى من ملاحقة فريق القطاع الخاص بالمشروع والممثل بشركة Celera Genomics، و التى كانت قد هددت بالاحتفاظ بالنتائج لنفسها لاستغلالها تجارياً. وبذلك تكللت الجهود المشتركة فى الإعلان عن نتائج المشروع الأولية بواسطة كلا الجانبين فى ٢٦ يوليو سنة ٢٠٠٠.

وبعيداً عن تعقيدات المعادلات الرياضية وبرامج الكمبيوتر فالعلمانية تبدأ أساساً من وجود قواعد للمعلومات أو بمعنى آخر من وجود نظام لتبويب وتنظيم المعلومات، ومن خلال هذا المفهوم فيمكن اعتبار المكتبات (الكتب والدوريات والمجلات والصحف) هى المكان حيث تجمع المطبوعات وتبويب وتنظم حسب معايير عديدة متباينة. بهذا المفهوم العلمانى، كان لمدينة الإسكندرية (Alexandria) شرف السبق فمنذ حوالى ٢٠٠٠ سنة مضت كانت الاسكندرية منارة للعلم ومركز ثقافى وتجارى متميز بجنوب شرق البحر الأبيض المتوسط، وكانت مكتبتها الشهيرة هى الأولى من حيث تنظيم وتبويب المعلومات والمخطوطات التى جمعت من كافة الحضارات مثل اليونانية و الأشورية و الفارسية والمصرية والعبرية والهندية وغيرها. فإذا ما كان لنا شرف السبق فوجب علينا الآن اللحاق بالدرب، فإلى هنا نضيّق ونلحق ما فاتنا.

وكان من النتائج المباشرة لتغلغل العلمانية فى العلوم البيولوجية أن ظهر فى السنوات العشر الماضية مجموعة من العلوم الجديدة والتى تتميز جميعها بالمقطع (-omics) أو ما يعرف مجازاً بعلوم "الأومكس" مثل الجينومكس (Genomics) والبروتيومكس (Proteimics) و الفينومكس (Phenomics) وغيرها. ويمكن تعريف الجينومكس بأنه الفرع من العلم الذى يهتم بدراسة المحتوى الوراثى ككل لكائن ما.

ومصطلح جينوم genome معروف فى الوراثة منذ أوائل الثلاثينات من القرن الماضى بين رجال الوراثة السيتولوجية cytogenetics، ويعنى العدد الأحادى من الكروموسومات haploid. وقد نوه العالم C. P. Blacker فى كتابه المعنون Eugenics

المنشور سنة ١٩٥٢ بأن ظهور مصطلحات وراثية مثل المجاميع الجينية والجينوم التي غيرت من المفاهيم الوراثية السابقة من حيث التركيز على الجين بصفته وحدة التوارث إلى الاهتمام بالوحدات الأكبر مثل الجينوم. ويمكن أيضا تعريف الجينومكس بأنه تقنيات "السلسلة" sequencing جزيئات الـ DNA الكلية (الجينومية) لكائن ما، ولكن إمكانية تطبيق تلك التقنيات لم تكن ممكنة قبل مجهودات عالم الوراثة العظيم Sanger سنة ١٩٧٧ حيث تمكن من سلسلة الجينوم الكامل للفيروس البكتيري الصغير *phiX174* والبالغ قدره ٥٣٨٦ زوج من النيكلوتيدات فقط، وقد حصل على جائزة نوبل فيما بعد على إنجازهِ هذا. ومع أن الطرق التي اتبعتها كانت بطيئة (بمفهومنا الحالي) إلا أنها كانت البداية التي فتحت الباب على مصرعية نحو بذوغ فجر الجينومكس، خصوصا بعد الإعلان عن استكمال سلسلة الجينوم الكامل للبكتيريا *Haemophilus influenzae* (١,٨٣٠,٠٠٠ زوج من النيكلوتيدات) سنة ١٩٩٥.

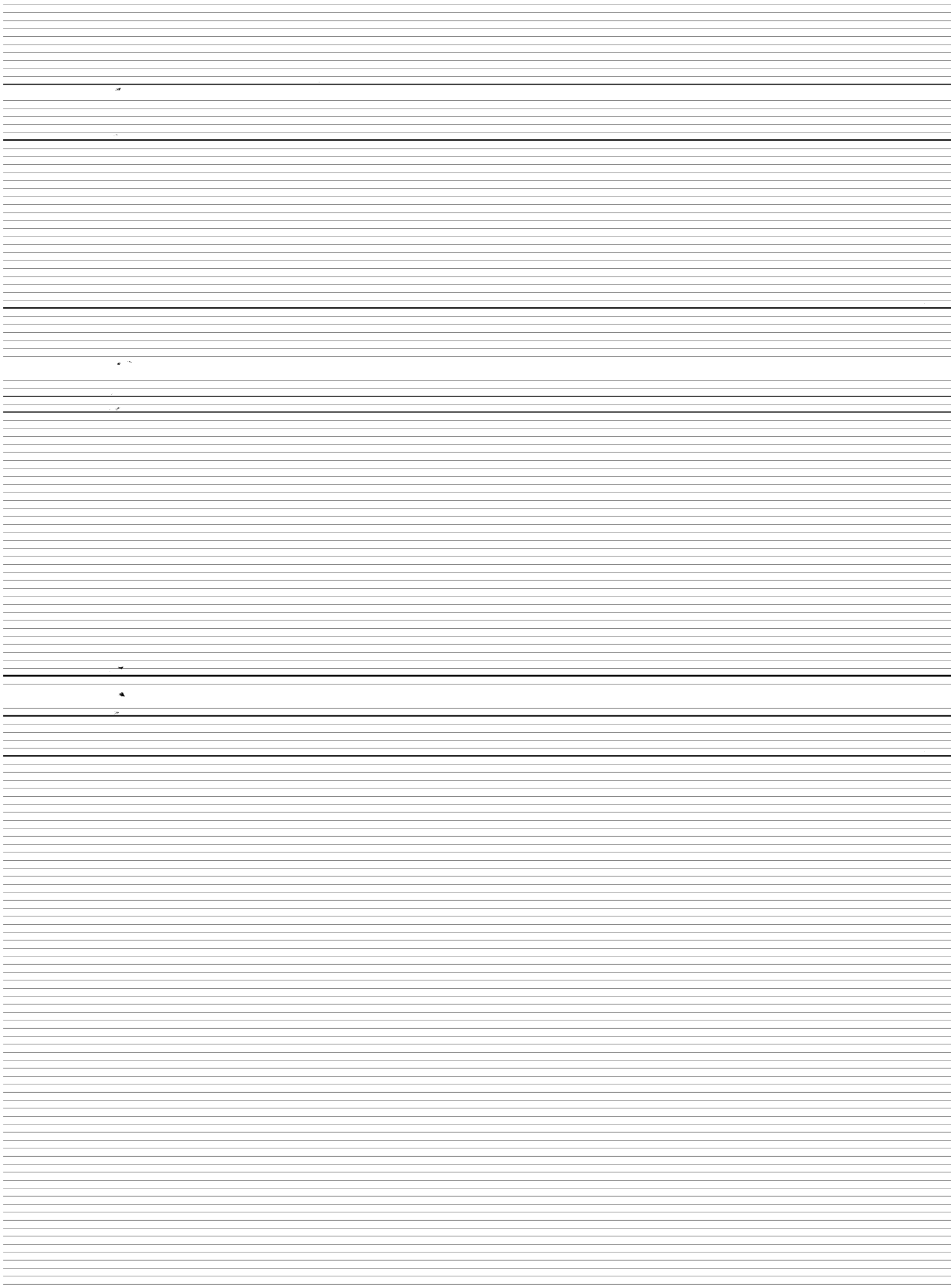
كذلك وجب التنويه، بأننا قد استعملنا المصطلحات الخاصة بهذا الموضوع دون محاولة تعريفها، بل هي المصطلحات الأجنبية مكتوبة بحروف عربية، ومع قناعتى بعدم صحة هذا الاتجاه إلا أن واقع الأمر فرضه علينا. فهذا المجال حديث للغاية، بل انه يعتمد أساسا على شبكة الأنترنت في الدراسة والتحصيل فكان لزاما أن يعرف الدارس المصطلح الذي يمكنه من الاتصال بتلك المصادر العالمية بسهولة. تاركين الباب في المستقبل، بعد أن ترسخ تلك المفاهيم الجديدة في ذهن القارئ العربي، لكي نعربها دون الإخلال بالمعنى ولكن تلك مهمة لاحقة قد تجد من يتصدى لها. ومع ذلك فقد حاولنا تعريف عدد من المصطلحات الوراثية والتطورية الجديدة فيما هو متاح، ولكن نستطيع القارئ عنرا في احتمال توارد أكثر من تعريف لنفس المصطلح لتعدد المشاركين في إنجاز هذا العمل.

ولا يفوتنا تقديم جليل الشكر للدكتور محمد عبد الفتاح ياقوت لعظيم مساعداته لإخراج وطباعة هذا العمل في أكمل وجه ممكن.

وختاما أتمنى أن يلقى هذا الجهد المتواضع اهتمام الدارسين وأن يكون سبيلهم للإضافة والبحث. والله من وراء القصد،،،

أحمد يوسف المتينى

الإسكندرية في ديسمبر ٢٠٠٥.



1

2

3

4

5

6

7

١. المعلومات الوراثية

Genetic Information

إعداد: أحمد المتينى

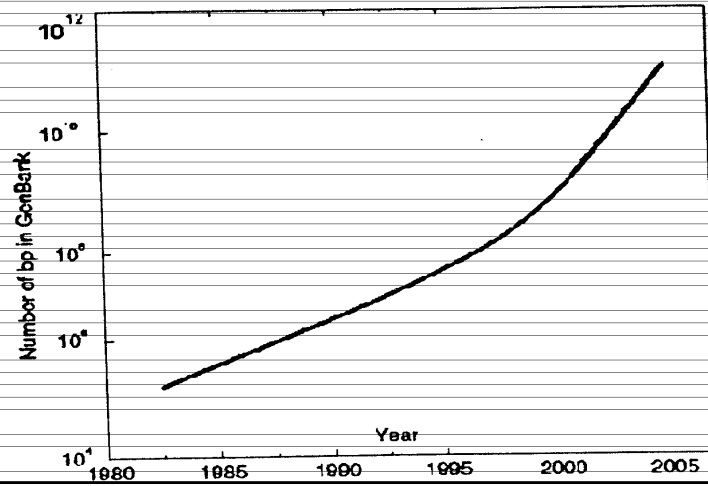
ركز علماء الوراثة، منذ أربعينيات القرن الماضى وحتى الآن، جل جهدهم وفكرهم فى التعرف على الطبيعة الجزيئية للجينات molecular nature of genes، حيث قتل الموضوع بحثا من كافة جوانبه تقريبا وتراكم لدينا كم هائل من النتائج أصبح الإناء بكافة جوانبها من شبه المستحيل حتى على المتخصصين. وزاد من تعقيد الوضع القائم التقدم الهائل خلال السنوات العشرة الأخيرة فى الانتهاء من دراسة جينومات العديد من الكائنات ومن ضمنها الإنسان، والجدول التالى يحدد عدد الجينومات التى تم معرفتها حتى

سنة ٢٠٠٢.

السنة	عدد الجينومات المدروسة لكل سنة	العدد التراكمي
١٩٩٤	-	-
١٩٩٥	٢	٢
١٩٩٦	٢	٤
١٩٩٧	٥	٩
١٩٩٨	٨	١٧
١٩٩٩	١٣	٣٠
٢٠٠٠	٢٢	٥٢
٢٠٠١	٤٢	٩٥
٢٠٠٢	١٠٠	١٩٥

وقد يبلغ عدد الكائنات التى تم التعرف على جينوماتها حتى نهاية عام ٢٠٠٥

ضعف هذا الرقم مما زاد بشكل مهول حجم البيانات المتاحة. شكل (١-١) يبين معدل الزيادة فى حجم البيانات المخزنة بالبنوك الجينية gene banks من عام ١٩٨٠ حتى نهاية عام ٢٠٠٥، مما يوضح بصورة جلية قدر المشكلة التى يعانى منها البيولوجيون عموما والوراثيون خصوصا. هذه الصعوبة دفعت العديد من العلماء والباحثين فى التفكير فى مخرج من هذه المشكلة. وكان الحل المأمول هو الاستعانة بعلوم الكمبيوتر والبرمجة.



شكل (١-١) : عدد أزواج النيوكليوتيدات المخزنة في البنوك الجينية خلال الخمسة وعشرون سنة الماضية.

الوراثة مثلها مثل أى طراز آخر من المعلومات information لا بد من أن تعتمد على لغة language أو وسيلة ما للتواصل ونقل المعلومة، لذلك فالمعلومات الوراثية لها لغتها الخاصة (ثلاثيات الشفرة الوراثية – genetic code) والتي لا تختلف كثيرا عن اللغات الحية (مثل الإنجليزية والفرنسية والعربية.... وغيرها) ولا تختلف أيضا عن لغة الكمبيوتر machine language كثيرا. والمعلومات الوراثية لا تختلف عن أى لغة أخرى، فهي تتميز بالصفات المعلوماتية العامة مثل:

- (١) أنها معلومات رقمية (digital signals) ، النيوكليوتيدات الأربع تكون ثلاثيات الشفرة (triplet-nucleotides - codons) المستقلة والتي ليس بينها تداخل (non-overlapping) حيث تمثل الحروف الهجائية لبناء البروتين
- (٢) المعلومات تمثل مصفوفات خطية (أوتار) (linear strings) مكونة غالبا من وحداتها الكودون المصفوفة والمتجاورة

(٢) تمتاز أيضا بوجود نظام للتحقق والتصفية (filtering) للتعرف على المعلومات ذات
المـدلول (sense-information) أو عديمـة المـدلول أو الشوشـرة
(non-sense or noise).

١- طبيعة الجينات Nature of Genes.

الوراثة هي العلم الذي يختص بدراسة إنتقال الصفات (المظهرية) من أفراد النوع
الواحد إلى أبنائهم من خلال التكاثر الجنسي. ومنذ مندل وحتى الآن فإن مفهومنا عن
طبيعة العامل الوراثي (الجين) المسئول عن نقل الصفات من الآباء إلى الأبناء، مبهم وغير
محدد! ففي القرنين السابع عشر والثامن عشر، ساد الاعتقاد بنظرية "سبق التكوين"
preformation theory حيث كانوا يعتقدون بتكون الفرد كاملا قبل المولد (أي عدم
الاعتراف بوجود الجين)، ولكن مع مطلع القرن التاسع عشر نادى عدد من العلماء وعلى
رأسهم دارون و ويزمان بوجود العامل الوراثي كوحدة مادية (particulate theory)
ولكن اعتقدوا بأن هذه العوامل يمكن المزج بينها، هنا جاء الدور العظيم لمندل حيث بين
أن هذه العوامل (الجينات) لا تمتزج ولكن تربطها علاقات من السيادة والتنحي. وبإعادة
اكتشاف تجارب مندل مع مطلع القرن العشرين، ساد المفهوم المندل للجين أو ما يسمى
بالجين الكلاسيكي (classical gene)، ولكن يجب أن نؤكد هنا أن المفهوم الكلاسيكي
للوراثة كان جل اهتمامه التنبؤ (prediction) بالأشكال المظهرية المتوقعة بفرضية وجود
الجينات المسئولة عنها دون معرفة طبيعتها (وهذا كان افتراضا نظريا بحتا)، أي الجين هي
بمثابة "صندوق أسود" حيث تمثل عوامل الآباء المدخلات inputs بينما الشكل الظاهري
يمثل المخرجات outputs، بينما محتوى هذا الصندوق فمجهول تماما، مع العلم أن
الرابطه بين الشكل الظاهري والجينات لابد أن تتم من خلال عمليات التشكل والنمو
(development) تبعا للفكر المندل ذاته. هذا التجاهل لطبيعة الجين المادية تداعا مع
منتصف القرن العشرين مع التوجه لدراسة البيولوجيا الجزيئية عموما والوراثة
الجزيئية خصوصا بفضل مجهودات عدد من علماء الوراثة العظام الذين ساهموا بدون
شك في تغيير الكثير من المفاهيم الوراثة التي سادت لأكثر من قرن من الزمان، بصورة
مباشرة أو غير مباشرة ويكفيهم فضلا أننا نعيش اليوم عصر الجينوم.

ولقد تراكمت المعارف خلال الخمسين سنة الماضية نحو ترسيخ مفهوم جديد للجين، فمع بدايات الخمسينات من القرن الماضي تم الإنجاز الأكبر في علم الوراثة، بل قل في علم البيولوجى عامة، وكشف الغطاء عن محتوى الصندوق الأسود (الا وهو الجين)، وكان الفضل يرجع لكل من واتسون و كريك سنة ١٩٥٢ في التعرف على التركيب الجزيئي للجين، وهو الأحماض النووية خصوصا الـ DNA. منذ هذا التاريخ حتى يومنا هذا تضاعفت معارفنا عن تركيب (structure) ووظيفة (function) الجين، فمن مفهوم جين واحد – إنزيم واحد إلى جين واحد - سلسلة ببتيدية واحدة - إلى مفهوم العقيدة الأساسية (central dogma) للتعبير الجيني التي تقر بأن الجين (gene) عبارة عن تتابعات من النيوكلووتيدات على طول الـ DNA تمثل إطارا (أو إطارات) للنسخ والترجمة (open reading frame) بالإضافة لعدة تتابعات عند النهايات ٢' و ٥' لاتنسخ ولكن تختص بالتحكم في التعبير الجيني. ويتم نسخ الجين إلى نسخة من الـ RNA المرسال، وبمساعدة جزيئات أخرى من الـ RNA تترجم هذه النسخة إلى سلسلة ببتيدية (بروتين) محدد له وظيفة معينة. ظل هذا المفهوم للجين صامدا لأكثر من نصف قرن، ولكن مع نهاية القرن العشرون وبداية القرن الحادي والعشرون تجمعت لدينا حقائق جديدة جعلتنا نتشكك في ثبات وصحة تلك العقيدة. ففي خلال السنوات القليلة الماضية، خصوصا مع مطلع القرن الحادي والعشرين، تجمعت أدلة وملاحظات تجريبية تجعل من صحة وشمولية هذه العقيدة الأساسية لمفهوم الوراثة مشكوك فيها، وفيما يلي بعض من تلك الحقائق التي دفعتنا إلى إعادة التفكير:

١. التحقق من أن عدد الجينات المشفرة للبروتين (ORFs) في الإنسان متدني عن المتوقع بشكل كبير (٢٠,٠٠٠ – ٤٠,٠٠٠ جين) بل تذهب بعض الدراسات الحديثة بعددها لأقل من ٢٥,٠٠٠ جين !! فكيف يفي الإنسان حاجته من البروتينات (حوالي ١٠٠,٠٠٠ بروتين) بهذا القدر الضئيل من الجينات، والذي قد يقارب ما هو موجود في كائنات أقل تطورا !! وأن قدر الجينات المشفرة لا تشكل أكثر من ١,٥ – ٢ ٪ من جملة الجينوم في الكائنات الراقية عموما.

٢. التعرف على عدد كبير (يزداد يوما بعد يوم) من الجينات غير المشفرة فيما يعرف بأسم ncRNAs، وهي التي تنسخ ولا تترجم إلى بروتين ولها العديد من الوظائف

التنظيمية بل والخلوية. وأتضح أنها تمثل حوالى ٩٧ – ٩٨ ٪ من حملة النشاط النسخي (transcriptomea) فى الكائنات الراقية. وإذا ما اعتبرنا هذه التتابعات (ncRNAs) جينات !! فإنها ستزيد قدر المعلومات الوراثية من حدود الـ ٢ ٪ إلى حوالى الـ ٤٠ - ٥٠ ٪ من حملة الـ DNA الخلوى، لكن سيتبقى قدر كبير من الجينوم (حوالى ٥٠ ٪) غير محدد الوظيفة، فيما يسمى بالنهايات (junk DNA) !! هل يمكن أن يكون هناك

نفايات فى النظام الحي ؟

٣. حالات الإعداد المفاير للمرسال الأولي (pre-mRNA alternative splicing) والتي تعددت صورها وأشكالها خلال الدراسات التي تمت خلال السنوات العشرين الماضية، وفيما يلي بعض منها:

(١) الإعداد المفاير من البداية (الطرف ٥') مثل جين (مجازاً) الألفا أميلاز فى الفأر أو الأعداد المفاير من النهاية (الطرف ٣') مثل جين السلسلة الثقيلة للجلوبين المناعي (immunoglobulin heavy chain) أو الأعداد المفاير من كلا الطرفين مثل جين تروبونين-ت فى العضلات (muscle troponin-T)، وكذلك حالات تحرير (تحويل) المرسال (mRNA editing) مثل جين apolipoprotein-B وغيرها.

(٢) الإعداد المتداخل لأطار واحد دون أدنى مشاركة بين المرسلات. Overlapping ORFs without shared coding sequences. وفى هذه الحالات ينسخ ويعد أحد الأطارات مشتملاً الأكسون الأول مثلاً فى الحالة الأولى ويستبعد فى الحالة الثانية لكونه يمثل أحد الأنترونات، ومثال هذه الحالات:

. DNA complex IP259/Dub80 in *Drosophila melanogaster*

(٣) الإعداد النسخى المشترك لأنئين من الإطارات المستقلة.

Co-transcriptional splicing between two ORFs مثال حالتي

P2Y11 & SSF1 فى الإنسان.

(٤) الإعداد النسخى المشترك المشتمل أحد الجينات الكاذبة.

. Co-transcriptional splicing with pseudogene

(٥) ومثال ذلك:

CYP3A7 coding sequence and the pseudogene CYP3AP1 on human chromosome 7q21-q22.1

(٦) الإعداد النسخي المتبادل المتوازي (لكلا ذراعي الـ DNA لنفس الإطار).

trans-splicing Alternative ومثالها:

The modifier gene (mdg4) in *D. melanogaster*

وغيرها وغيرها من الحالات التي ينسخ ويعد فيها الإطار الواحد ليترجم لأكثر من بروتين (واللائحة كل يوم في إزدیاد)، ولكن أكثر الأمثلة الصارخة في هذا الصدد فهي جينات العائلة الجينية المعروفة باسم N-CAM حيث تكون مسئولة عن بناء طرز مختلفة من البروتين Neural Cell Adhesive Molecules، وهي أساسا مسئولة عن تكوين الذاكرة طويلة الأمد في مخ الإنسان إلا أنها توجد في كافة الأنسجة للتعرف والربط بين خلايا النسيج المتناظرة. المهم وجد أن أحد هذه الجينات يمكن أن يعد بطرق مختلفة قد ينتج عنها حوالي ١٠٠ طراز مختلف من البروتين iso-form }}

هذه الحقائق الجديدة لم تتحدى مفهوم "الجين الجزيئي" وحسب بل أسقطت الاعتقاد السائد منذ النُدلية، بأن الجين وحدة مستقلة محددة لها ذاتيتها الخاصة، فقد اتضح أن الجين قد يشارك على المستوى الجزيئي في أكثر من وظيفة واحدة، تركيبية أو تنظيمية. أن الهزة الكبيرة التي أصابت مفهوم الجين الجزيئي، والذي كنا نعتقد أنه صامد إلى الأبد، جعلتنا نعيد التفكير في مصداقية المسار الذي نحن سائرون فيه (التبحر في الدراسات الجزيئية عن الجينات). فمع تسليمنا المطلق بالأساس الجزيئي للجينات، إلا أن الواقع يؤكد أن هذه الجينات في نهاية المطاف هي أدوات لتخزين وتناول وتبادل "المعلومات information" وأن لغة التفاهم هي الشفرة الوراثية (genetic code) مثلها مثل أي لغة حية (مثل الإنجليزية أو العربية..... أو غيرها). وبناء على هذا التوجه فيمكننا أن نعرف الجينات بأنها: وحدات معرفية — تحمل المعلومة لصفة ما ولكنها ليست منطقة تلقائيا بتنفيذ تلك المعلومة.

٢.١. المعلومات من تتابعات الـ DNA. DNA Sequence Information.

من المعروف لنا أن جزيئات الـ DNA هي الحامل للمعلومات الوراثية في الغالبية العظمى من الكائنات (الاستثناء في ذلك عدد قليل من الفيروسات حيث يحل محلها جزيئات من الـ RNA)، فهل تلك المعلومات محصورة فقط فيما نعرفه من معلومات العقيدة الوراثية الأساسية (الجينات المشفرة coding sequences) ؟ أم هناك معلومات أخرى خلافاً لذلك ؟ لقد اتضح لنا أن هذه الجزيئات تحمل قدر من المعلومات أكبر مما كنا نتوقع، عرفنا بعضها ومازال أغلبها مجهولاً. وفيما يلي ملخص لأهم المعلومات البيولوجية التي يمكن للـ DNA القيام بها أو التحكم بها:

- التتابعات المشفرة لأحماض أمينية (بروتين) وذلك إما :
 - بطريقة مباشرة مثل أغلب الجينات المعروفة mRNAs.
 - بطريقة غير مباشرة مثل حالات:
 - التجمع للجينات في البروتوزوا- تغير على مستوى الـ DNA.
 - التحرير على مستوى الـ RNA editing . RNA.
 - تفصيص الـ RNA splicing . RNA.
 - تفصيص البروتينى Protein splicing .

(الإنسولين).

- التتابعات غير المشفرة non-coding ، ومنها :

- الناقل tRNA.
- الريبوزى rRNA.
- النووى الصغير snRNA.
- التيلوميرازى telomeraseRNA.
- وغيرها مثل RNAi

- التتابعات لبروتينات الارتباط بمواقع محددة على الـ DNA

protein binding sites.

- التتابعات الخاصة لتحديد الشكل الفراغي للجزئ، ومنها:
 - الألتفافات في الحلزون intrinsic helix curvature.
 - تحديد مواقع النيوكلوسومات nucleosome positioning.

- تتابعات تحدد الثبات التركيبي الوظيفي، ومنها:
 - تحديد بداية النسخ transcription initiation.
 - نقطة أصل التضاعف origin of replication.
 - أماكن التطفر الساخنة "hot spots" mutational.

إذا ما كانت هذه بعض من الوظائف التي قد يقوم بها الـ DNA، فهل هناك اختلافات تركيبية لهذه الجزيئات قد يكون لها إنعكاسات وظيفية (قد لا نعلمها الآن ولكن سنعلمها في المستقبل) ؟ وهل هذه الاختلافات التركيبية تعكس معلومات ما ؟ مثل هذا التساؤل يمكن فحصه من خلال عمل مسح scanning لتتابعات الـ DNA بحثاً عن مناطق متشابهة مميزة بذاتها أو ما نطلق عليه التكررات المتجاورة repeats. كذلك البحث عن أشكال مغايرة للتوزيعات الفراغية لجزيئات الـ DNA، أو أشكال الحلزون المختلفة الناجمة عن متكررات متجاورة من البيورينات أو البيريميدينات.

١.٢.١. التكررات Repeats.

من دراسة تتابعات الـ DNA أتضح لنا وجود معدلات عالية من التكررات على طول هذه الجزيئات يمكن تقسيمها إلى الآتي:

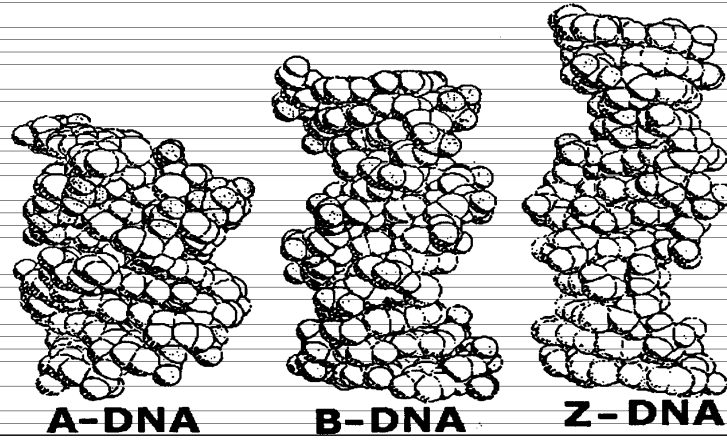
- متكررات مباشرة Direct repeats $\Rightarrow \Rightarrow$ ، يوجد منها:
 - متكررات صغيرة متجاورة simple tandem repeats.
 - متكررات أكبر متجاورة longer tandem repeats.
 - متكررات غير متجاورة non-tandem repeats.
- متكررات متواجهة Phased repeats \Rightarrow ، يوجد منها:
 - متكررات منعكسة inverted repeats \Rightarrow .
 - متكررات صورة في المرآة mirror repeats $\Rightarrow \Leftarrow$.
 - متكررات أرتدت reverted repeats $\Rightarrow \Rightarrow$.

وكل هذه الطرز من التكررات repeats تنتشر في جينومات كافة الكائنات التي درست حتى الآن بصورة كبيرة خصوصا في مناطق السنتروميرات centromeres والتيلوميرات telomeres الكروموسومية. وتتواجد بعض من صورها أثناء عمليات العبور الوراثة recombination (في المناطق ثلاثية الأذرع - triplex) وكذلك أثناء التضاعف replication. هذه التراكيب قد يكون لها فائدة عند إجراء المقارنات بين الجينومات أو التتابعات المختلفة.

٢.٢.١ نماذج الحلزون DNA helix models.

لوحظ وجود أشكال مختلفة لحزون الـ DNA، ولكن أكثرها وجودا وتكرارا في الكائنات الحية هو النموذج التقليدي الذي اقترحه واتسون وكريك سنة ١٩٥٢ والمعروف لنا جميعا، لذلك سمي فيما بعد بالنموذج البيولوجي biological-DNA ويعرف اختصارا بأسم B-DNA. بالإضافة للنموذج البيولوجي وجدت نماذج أخرى مختلفة وتوجد في الأنظمة الحية ولكن بندرة عرف هما بأسم نموذج A-DNA ونموذج Z-DNA. وشكل (٢-١) يوضح هذه النماذج الثلاث كما نشرها Dickerson et al. سنة ١٩٨٢ لقطع متساوية لكل منها.

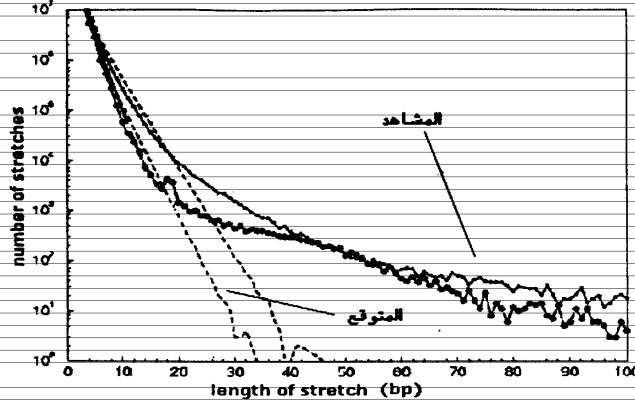
عائلة الحلزونات من الطراز A-DNA أكثر تواجدا في جزيئات الـ RNA مزدوج الذراع ds-RNA أو الهجن بين الـ DNA والـ RNA (RNA/DNA hybrids). ويمكن أن يتواجد من خلال بعض التتابعات حيث تتكرر تتابعات متجاورة (خمس نيوكلووتيدات) من البيورينات purines. أما عائلة الـ B-DNA فهي الغالبة في أنوية الكائنات الحية، ولكنها قد تظهر قدرا من الاختلافات، فمن المعروف أن ١٠ أزواج من النيوكلووتيدات المقترنة تمثل لفة واحدة للحزون المزدوج، لكن وجدت تتابعات محدد قد تؤدي لتكوين اللفة مابين ٩ و ١٢ نيوكلووتيدة مقترنة. العائلة الثالثة هي عائلة Z-DNA وهي أقل الثلاث حدوثا، ولكن عند تتابعات معينة مثل متكررات GC (GCGCGC)



شكل (٢-١): رسم توضيحي لنماذج الـ DNA (A-DNA, B-DNA, Z-DNA)، كل ممثل بقطعة طولها ١٢ زوج من النيوكليوتيدات المقترنة.

يتكون الحلزون Z-DNA عادة في الكائنات حقيقية النواة eukaryotes. كذلك لوحظ أن بعض الجزر من CpG العادية أو المضاف لها مجاميع المثل تكون هذا الحلزون في التجارب العملية، كذلك عزلت بعض البروتينات التي تفضل الارتباط بمناطق الحلزون Z-DNA.

الملفت للنظر أن الدراسات الحديثة أوضحت أن حدوث وتكرارات هذه الحالات في البكتيريا اعتباطي ويرجع للصلفة ولكن وجد أن تكرار هذه الحالات غير اعتباطي وينحرف معنوياً عن المتوقع في دراسة على جينوم الإنسان، كما هو موضح بشكل (٢-١).



شكل (٢-١): علاقة بيانية بين تكرارات حالات الحزوز الثلاثة وأطوالها كما هو متوقع وكما هو مشاهد. لتتابعات كروموسوم ١ في الإنسان.

حقيقة كون أن تكرار هذه الحالات في الكائنات حقيقية النواة لا يرجع للصدفة، بينما الحال مختلف في الكائنات غير حقيقية النواة حيث وجد أن أغلبها يرجع للصدفة، هذه النتائج تلفت النظر إلى أهمية الصدفة في دراستنا للأنظمة البيولوجية. لذلك وجب التعرف على مفهوم الصدفة ولم بصورة مبسطة.

١. ٢. ٢. التركيب الثانوي لجزيئات الـ RNA.

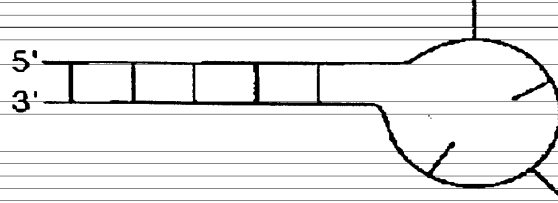
RNA Secondary Structure.

جزيئات الـ RNA على اختلاف أشكالها وطرزها تختلف عن جزيئات الـ DNA وإن كانت لا تقل عنها أهمية من حيث دورها في نقل وتداول المعلومات الوراثية. جزيئات الـ RNA تختلف عن الـ DNA في كونها جزيئات خطية (مفردة النراع) من تجمع وحداتها.

الأقترانات التقليدية ما بين الجوانين والسيتوسين والأدينين واليوراسيل، ولكن هناك أيضاً اقترانات غير تقليدية أكثرها شيوعاً ما بين الجوانين واليوراسيل. هذه الأوضاع والأشكال الجزيئية الناجمة عن اقترانات تؤدي إلى ما يسمى بالتركيب الثانوي secondary structure وهو حالة وسطية بين التركيب الخطي والتركيب ثلاثي الأبعاد.

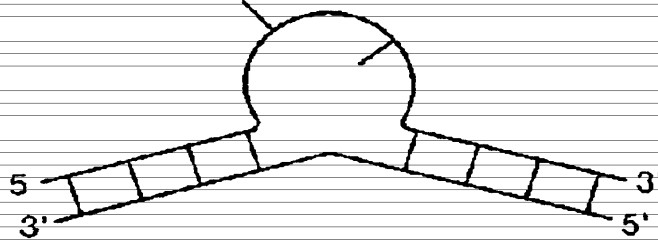
3D structure. وبناء على التركيب الثانوي تتكون عدة تراكيب جزيئية مميزة يمكن إجمالها في الحالات التالية:

- الانتفاخ القاعدي stem loop أو انتفاخ دبوس الشعر hairpin loop، كما هو مبين بالرسم التالي:

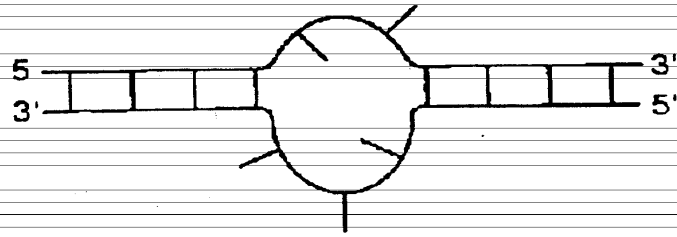


وعادة لا يقل الانتفاخ عن أربعة نيوكلووتيدات غير مقترنة.

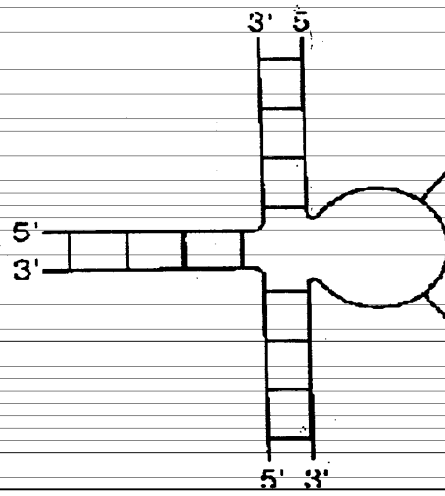
- الانتفاخ البرميلي Bulge loop، ويتكون عادة عند مناطق عدم الإقتران في طرف واحد من الإنتفاخ، كما هو مبين بالرسم التالي:



- إنتفاخ داخلي interior loop يتكون عند مناطق عدم الإقتران على طرفي الإنتفاخ، كما هو مبين كالآتي:



- إنتفاخ الوصلة junction loop، ويتكون من اقتران جزيئين من الـ RNA مكونان تراكيب ثانوية مشتركة، كما هو مبين كالتالي:

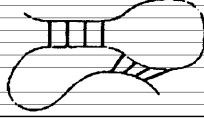


- التراكيب الثلاثية tertiary structures وهي تنتج من الاندماج بين أكثر من جزئ، ومنها اشكال عدة يمكن تمثيل أهمها كالتالي:

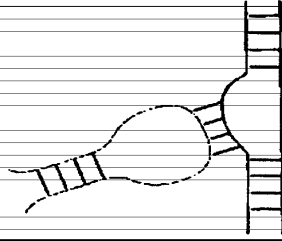
- إنتفاخات متقابلة kissing loops



- العقد الكائبة psuedoknots



- التداخل بين إنتفاخ برميلي وديوس الشعر hairpin and bulge interaction



٣.١. المعلومات والصدفة و التحكم. Information & Chance & Control.

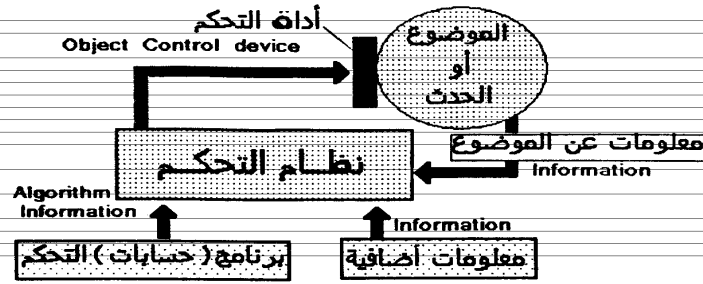
نحن نتعايش في كل لحظة من حياتنا مع عدد من الأحداث غير المتوقعة (unexpected) أو غير المتحكم فيها (uncontrolled) والتي يمكن أن تندرج تحت ما نعرفها عليه بأسم الصدفة (chance)، وهي تحدث دوما في العالم من حولنا فهي جزء لا يتجزأ من ظواهر الكون. فهل تكرار الصدف من حولنا دفعنا يوما للتفكير في طبيعة

هذه الصدفة ؟ في أغلب الأحيان لانتلفت كثيرا لتلك الصدفة من حولنا، وهذا ليس نوعا من الإنكار لوجودها ولكن لتألفنا معها واعتبارها أمر مفروغ منه. ولكن وبدون أدنى شك لقد حاول البشر معرفة لماذا يحدث غير المتوقع من الصدفة، بل يمكننا القول بأن الحضارة الإنسانية منذ فجر التاريخ وحتى الآن هي في سعى دؤوب لمعرفة غير المتوقع (الصدفة) أو المجهول. فيمكننا القول بأن الصدفة هي مرادف عدم المعرفة أو الجهل، فكلما زادت معلوماتنا كلما تضائلت فرص الصدفة من حولنا. لتوضيح الأمر أكثر دعنا نسوق المثال التالي : إذا ما كنت في الخارج راجعا إلى منزلك ظهرا بعد دوام العمل كما هي عادتك يوميا – وضغطت على جرس الباب الخارجى لمنزلك – فإنك عادة تكون متأكدا بقدر كبير بأن أحدا (زوجتك مثلا) سيفتح لك الباب، ولكن في يوم من الأيام لم يفتح لك أحد قبل بعد عشر دقائق، هذا الأمر غير المتوقع (الصدفة) سيصيبك بالدهشة بل بالإنزعاج لأن أحدا ليس بالمنزل – ولكن إذا ما شككت في إنقطاع التيار الكهربى عن شقتك – وبدأت بالطرق بيدك ولم يفتح لك أحد – سيزداد إنزعاجك من المصادفة خصوصا إذا اتضح لك أن نور السلم يعمل وتزداد حيرتك من وقع المصادفة (كل ذلك راجع لجهلك بالأمور الجارية). في أثناء ذلك هداك تفكيرك أن تتصل بزوجتك على "المحمول" فتخبرك بأنها اضطرت لترك المنزل لطارئ ما – هنا تطمئن وينتهى تأثير الحادثة عليك وتتقبل الأمر ببساطة. إذن انتهى تأثير الصدفة (الأحداث غير المتوقعة) عليك بمجرد معرفتك بحقيقة الأمر، أى أن المعرفة تقلل من احتمال الصدفة من حولنا. كذلك مع توفر وسيلة للتحكم control قللت من احتمال حدوث الصدفة. ففى مثالنا السابق – صادف هذا الرجل حدث غير متوقع (صدفة) كادت أن تخرجه عن إترانه، ولكنه تغلب عليها، فلماذا ؟ لأنه حاول التحكم فى الحدث مستعملا عقله كأداة للتحكم، ثم إستعان بوسائل متاحة لزيادة معلوماته (الاتصال بالمحمول أو إنارة نور السلم وغيرها....). دعنا الآن نخرج من حدود هذا المثال الضيق، لاستخلاص بعض القواعد العامة، فتبعا لقانون الديناميكا الحرارية الثانى: ففى الأنظمة المغلقة، فإن النظام يصل إلى حالة من استنزاف للطاقة تسمى "الأنتروبيا" (entropy) بعدها ينعدم أو يفنى النظام. أما نظرية التطور (على الجانب الأخر) فتقر بتعقيد النظم الحيوية وانتظامها الدقيق بل الفائق وتسلسلها عبر الزمن من منشأ واحد. هذا التعارض البين بين كلا النظريتان حول

مصير الأنظمة الحية، من الفناء المحتوم تبعاً لقوانين الديناميكا الحرارية - أو البقاء والتعقيد المتحكم فيه تبعاً لنظرية التطور. هذا التعارض دفعنا للتفكير لاستخلص من تلك الظواهر والحقائق المادية ، ملاحظة غاية في الأهمية وهي: "أن التحكم (control) وتوفير المعلومات (information) في الأنظمة عامة (والبيولوجية خاصة) يؤدي إلى تقليل الأنثروبيا ويسمح بقدر من الإنتظام (order)". إذن التحكم والمعرفة هما وسيلتنا لتقليل الصدف. شكل (٤-١) يلخص وسائل تقليل الصدف في الأنظمة المختلفة (مادية كانت أو حيوية). فإذا ما تعاملنا مع حدث (موضوع) ما بفرض تقليل الأنثروبيا (تقليل الصدف)

فلا بد من وجود أداة للتحكم (مثل الثرموستات في جهاز التكييف يعمل عند درجة حرارة محددة)، لكن لا بد من وجود عقل للتحكم قد يمثل برنامج حسابي algorithm يتحكم في تلك الأداة. لكن لكي يعمل هذا البرنامج لابد من مده بالمعلومات، أولاً عن الفرض من التحكم (ضبط درجة حرارة الغرفة) وكذلك معلومات إضافية (عن ظروف الجو خارج الغرفة مثلاً وغيرها وغيرها). إذا توفر ذلك فستكون لنا فرصة كبيرة للتغلب عن الصدف غير المتوقعة (تقلب الطقس مثلاً).

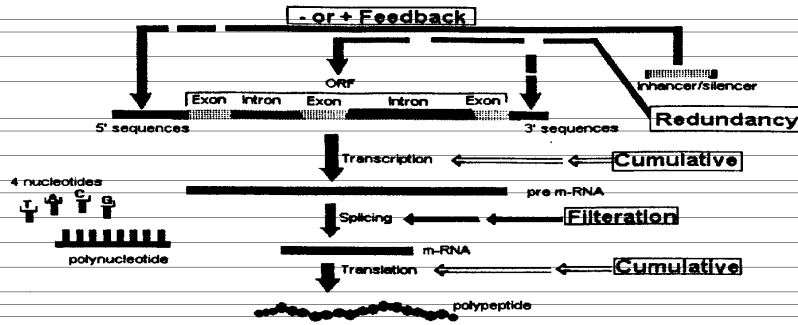
مما سبق يتضح لنا أن حرب البشر المستمرة على غير المتوقع من الصدف (أو عدم الإنتظام - disorder) كان وما زال لها جبهتان رئيسيتان للقتال، الأولى مباشرة للقضاء على الصدف (الشوشرة) العشوائية وقد أثبت التاريخ استحالة نجاح هذه الجبهة، أما الجبهة الثانية فهي الدبلوماسية (إذا جاز التعبير) بمعنى التعايش مع الصدف (التداخل العشوائي) ومحاولة التقليل من تأثيراتها دون إنكارها أو محاولة التقليل من قدرها، ويمكن تمثيل الاستراتيجية الثانية بحالة الحادثة في الهاتف في ظل وجود تداخل بالخط (شوشرة) فلا يمكننا التخلص منها أثناء الحادثة ولكن نتغلب على الوضع يمكن أن نرفع أصواتنا أو نكرر ما نقول عدة مرات، وقد أثبتت الأيام أن الجبهة الثانية هي المتاحة لنا في حربنا على غير المتوقع من الأحداث.



شكل (٤-١) : رسم تخطيطي عن وسائل تقليل غير المتوقع (الصدفة) لحدث ما.

ففى حربنا الدبلوماسية على الصدفة فإن وسائلنا عديدة مثل التكرار (redundancy) للمعلومات و تراكمها (cumulation) وغربلة (filtration) الصالح منها مع توفر أنظمة للتحكم (control). وهذه الوسائل نجدها ممثلة وبوضوح فى أى نظام وراثى، كما هو مبين بشكل (٥-١).

ففى هذا النظام توجد أنظمة متعددة موجبة أو سالبة للتحكم (سبق دراسة عدد كبير منها فى مقررات وراثية أخرى، ولا يسع المجال هنا لتكرارها). وكذلك نلاحظ التكرار للمعلومة الوراثية الواحدة (هناك نسخ متكررة من الجين الواحد)، كذلك فإن تراكم وتجميع المعلومة الوراثية عدة مرات ظاهرة عامة، ينسخ الجين الواحد ويترجم عدة مرات قد تصل مئات المرات للتأكد من وصول المعلومة المطلوبة. وظاهرة الانتقاء للأكسونات دون الأنثرونات خلال عمليات الأعداد والتفصيل للمرسل الأولى ماهى إلا وسيلة للغربلة والتصفية للتخلص من غير المرغوب من المعلومات مع الأبقاء على المرغوب منها.

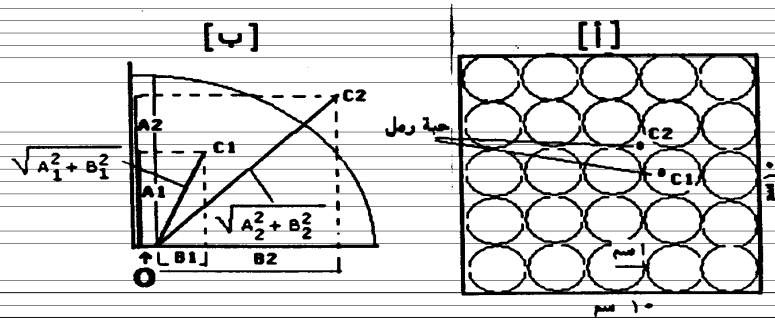


شكل (٥-١) : رسم توضيحي لنظام جيني عام يوضح نقاط التحكم وتقليل الشوشرة المتاحة.

مما سبق وضع لنا أهمية التعايش والقبول بوجود الصلابة في كافة مناحي الحياة من حولنا، ثم وجب علينا بعد ذلك محاولة التقليل من تأثيراتها ما أمكن. في محاولتنا هذه نحل اشبه بمن يلعب لعبة من لعب التسلية (a game) فعلى أن نحدد خياراتنا بين أفعال ولا تفعل وعلى إختيارى هذا ستتوقف النتيجة بين الفوز والخسارة، ومن هنا ظهرت "نظرية اللعبة" (Game theory) التى أصبح لها تطبيقات واسعة خصوصاً فى علوم الاقتصاد والأجتماع والسياسة. إذن علينا حساب القرار الصائب، وهذه الطريقة (وغيرها من الطرق الحسابية والأحصائية) عليها أن تزيد من احتمالات النجاح وتقليل احتمالات الفشل عن طريق حساب الاحتمال (probability) عند مستوى مصداقية محدد. فنحن فى حياتنا نتعامل من الاحتمالات (الحلم المقبول) بصورة مستمرة دون أن نلتفت لذلك، فعلى سبيل المثال - إذا ما سألت عن مساحة الدائرة - فستكون أجابتك المباشرة مساحة الدائرة تساوى مربع نصف القطر مضروباً فى π (πr^2) فإذا كان لديك دائرة نصف قطرها - ١ سم، فإن مساحتها $1 \times 3.14 = 3.14$ سم^٢. فمساحة الدائرة تقريبية إذن حيث أن π هنا قيمة مقربة، فهل خطر فى ذهنك يوم هل يمكن أن احسب قيمة π الحقيقية ؟ وهل يمكننى ذلك ؟ هنا يأتى دور طرق حساب الاحتمالات الحسابية، ولتبسيط انكثرة أكثر، دعنا نستعمل أحد الطرق الشهيرة وهى طريقة (Monte Carlo method) نسبياً للمدينة الشهيرة بنوادى القمار وخصوصاً لعبة الروليت التى بنيت عليها هذه الطريقة.

ولكن قبل ذلك دعنا نحل المسألة تجريبيا، فلو أحضرنا قطعة من ورق الكارتون أبعادها 10×10 سم أى مساحتها 100 سم² ورسمنا عليها دوائر كل نصف قطرها اسم بحيث تكون متماسة تماما (كما هو موضح فى شكل ٦-١) فسيكون عدده ٢٥ دائرة، ثم قطينا بحبة رمل صغيرة والقيناها على الرقعة بعدد كبير من المرات، ثم حسبنا عدد المرات التى وقعت فيها حبة الرمل داخل دائرة من الدوائر (C_1) وعدد المرات التى وقعت فيها بالمسافات البينية بين الدوائر (C_2)، تجدر الملاحظة أننا فى هذه التجربة تمثّلنا لعبة الروليت الشهيرة. فقل كانت نتيجةلقاء حبة الرمل ١٠٠٠ مرة - ٧٠٠ وقعت داخل أحد الدوائر و٣٠٠ مرة فى المسافات البينية، فيمكن القول أن مساحة الدوائر بالرقعة تمثل $100/700 = 10/7$ ، أى يمكن حساب مساحة الدوائر $100 \times 10/7 = 70$ سم² وحيث أن عدد الدوائر الكلى = ٢٥ إذن يمكن حساب مساحة الدائرة الواحدة ومنها يمكن حساب قيمة π ، والتى ستساوى فى هذه الحالة

٢,٨



شكل (٦-١) : محاولات لحساب قيمة π (١) تجريبيا و (ب) حسابيا.

اتضح لنا مما سبق أنه يمكننا أن نحسب تجريبيا قيمة π حسب الاحتمال المحسوب، فهل يمكن وضع نموذج حسابي (mathematical model) لتقدير الاحتمال دون اللجوء للتجربة خصوصا أن مثل هذه التجارب عرضة للأخطاء التجريبية العديدة (فعلى سبيل المثال عدم التأكد من مكان حبة الرمل إذا ما وقعت على الحدود وكذلك احتمال عدم إستواء الرقعة بنفس الدرجة وغيرها من العوامل غير المتحكم فيها). لنعمل هذا

النموذج الحسابى افترض اختيار رقمين عشوائيا مثل A و B بحيث ان يكونا اصغر من الواحد الصحيح واكبر من الصفر ($0 < N < 1$) و للتحقق من وقوعهما داخل او خارج احد الدوائر دعنا نتصور ربع احد تلك الدوائر، كما هو موضح بشكل (١-٦). فمن النقطة O اختار رقمين عشوائيا A_1 و B_1 وحدد نقطة تلاقيهما C_1 . اعد المحاولة وحدد النقاط A_2 و B_2 و C_2 . فالنقطتان C_1 و C_2 مواقع محتملة لحبة الرمل فى التجربة السابقة وتبعا للعلاقات الحسابية يمكن التوصل للمعادلتين التاليتين، إذا ما كان $A^2 + B^2 \leq 1$ فإن الرقم سيقع داخل الدائرة لكن لو كان $A^2 + B^2 > 1$ فإن الرقم سيقع خارج الدائرة. اذن يمكننا الآن حساب احتمال كلا الحدين دون رقعة ولا حبة الرمل ولا التكرار فى رمي الحبة. بدون شك ان النموذج الحسابى افضل، فبه نتخلص من الكثير من الأخطاء التجريبية التى تؤثر على صحة التقدير، ولكن مازال علينا ان نختار تلك الأرقام العشوائية وحساب المعادلات لكل منها مئات المرات مما يضعنا امام عمل شاق قد تفوق صعوبته رمي حبة الرمل العديد من المرات، أى لم نوفر شيئا فى الجهود والوقت. هنا جاء دور الكمبيوتر فهو يمكن ان يولد عدد لامتناهى من الأرقام العشوائية وحساب نتائج المعادلات فى دقائق بل فى ثوانى فى بعض الموديلات الحديثة. الاعتماد على الكمبيوتر اذن ضرورى ولكن الكمبيوتر لا يمكنه القيام بتلك المهام من تلقاء نفسه بل يجب كتابة تلك الأوامر والحسابات فى صورة برنامج رياضى يوضح للكمبيوتر مهامه المطلوبة منه، هذا البرنامج لحساب احتمالات الأحداث هو ما تعارفنا عليه بأسم طريقة مونت كارلو فى المثال السابق.

٤.١. الشبكات الجينية Gene networks.

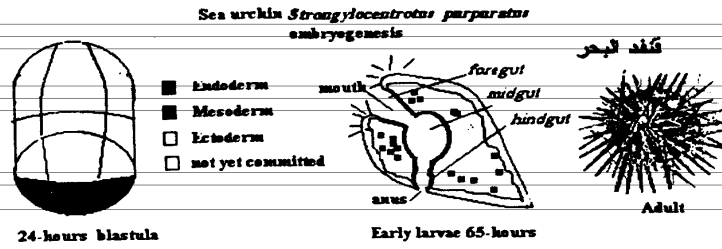
لقد عايشنا خلال أواخر القرن الماضى وأوائل القرن الحالى بذوغ عصر الجينومكس والمعلوماتية الحيوية وإنجازاتها الموهلة فى مجال تعرفنا على طبيعة المعلومات الوراثية، ولكن أدى بنا هذا الحال (زيادة المعلومات وتعقدها) إلى التخبط فى فهمنا للمنظومات الجينومية وكيفية تفاعلها وتداخلها مع بعضها البعض للوصول للشكل المظهرى. هذا الوضع دفع عدد من علماء الوراثة للقول بأن القرن الحادى والعشرين هو عصر ما بعد الجينومكس post genomic area، وفى تعليق لأحد أشهر الوراثيين وهو Lee Hartwell

(الحائز على جائزة نوبل سنة ٢٠٠١ لأكتشافه بروتينات الـ cyclins ودورها في دورة إنقسام الخلية) قال أن طلاب الوراثة في المستقبل سيكون مهمهم التعرف على الدوائر والشبكات الجينية gene circuits & networks عوضا عن دراسة الجينات بذاتها، ويجب علينا أن نتجه نحو بناء النماذج الوراثة أو ما يسمى بأنظمة المحاكاة simulation - systems في مجال الوراثة، فلو توفر لنا نظام يحاكي الخطوات الوراثة للأصابة بالسرطان مثلا، فسنكون أكثر قدرة على تفهم هذا المرض بل سنكون قادرين على تصميم العقاقير المناسبة لكل حالة بحالتها.

إن محاولات الباحثين لبناء الشبكات الجينية لحالات بعينها مازالت في مراحلها الأولى بل أن الإنجازات المتاحة في هذا المجال نادرة ومازالت في مراحلها الأولى وبعيده عن الوضوح والاكتمال، ولكن في سبيل بناء تلك الشبكات والدوائر الجينية يجب أن لا ننسى أنها نماذج نظرية وتختلف عن الأنظمة الحية، وهدفنا من بناء تلك النماذج سيكون فقط لزيادة قدرتنا على التنبؤ بالخرجات outputs المتوقعة من تلك المدخلات inputs الجينية.. وفي محاولتنا لتفهم هذا الموضوع سنكتفى بمناقشة أحد الدوائر لجين واحد فقط، لتجنب التعقيدات الموهلة التي قد تقابلنا في هذا المجال.

١.٤.١. الجين endo16 ومراحل تكوين جنين قنفذ البحر.

أختص عالم الوراثة الشهير Eric Davidsons وتلاميذه ومساعديه في معهد كاليفورنيا للتكنولوجيا بدراسة هذا الجين من أوائل تسعينات القرن الماضي وحتى الآن، حيث أجروا آلاف التجارب لعشرات السنين للوصول لتفهم هذا الجين، والذي يعتبر أكثر الجينات دراسة حتى الآن. قنفذ البحر (الريـتـزا) *Strongylocentrotus purpuratus*, sea urchin معروف في تجارب التشكل والتكوين، فمراحل تشكله الجنيني معروفة بالتفصيل (الأجنة شفافة يسهل تتبعها ميكروسكوبيا). فبعد الإخصاب تنقسم الخلايا الجنينية مكونة عدة طبقات تعرف بأسم ectoderm, mesoderm, endoderm وتشكل مرحلة البلاستيولا blastula من حوالي ٥٠٠ خلية بعد حوالي ٢٤ ساعة، بعد هذه المرحلة يتبعها مرحلة الجاستريولا gastrula حيث يتم إنتظام خلايا طبقة الـ endoderm للداخل لتكون جدران الأحشاء الداخلية لليرقة وذلك في حوالي ٦٥ ساعة لتكون اليرقة، كما هو موضح بشكل (١-٧).

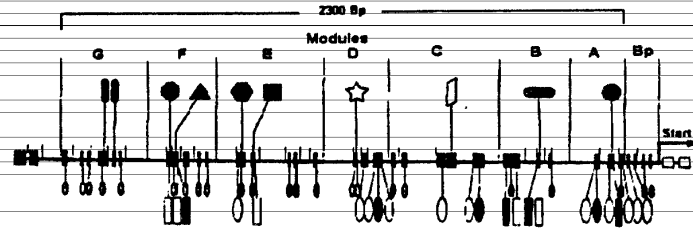


شكل (٧-١) : مراحل التكون الجنيني في حيوان قنفذ البحر *S. purpuratus*.

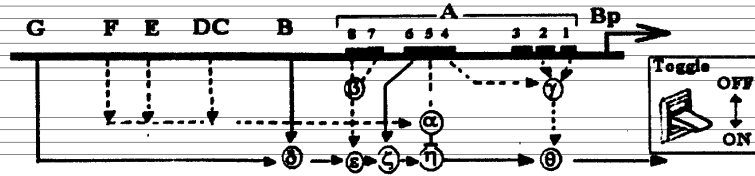
ومن دراسات Davidson *et al.* أمكن التعرف على الجين endo16 وهو الجين المسئول خلال مراحل الـ gastrula على تشكل خلايا الـ endoderm لتغليف الأحشاء الداخلية وأن نشاطه يحدد زمنياً ومكانياً في هذا النوع من الخلايا دون غيرها لتشكل منظومة تحكمية غاية في الدقة والتحكم. فقد وجد أن مناطق البروموتور البعيدة (المزيد من التفصيل يمكن الرجوع إلى المتينى - ١٩٩٧) تمتد على طول المناطق السابقة upstream لهذا الجين بطول حوالى ٢٠٠٠ زوج من القواعد المقترنة في المناطق البين جينية intergenic أو ما كان يعرف بأسم النفاية junk DNA (بدون وجه حق) وهي مقسمة إلى ثمانية مناطق متتالية من الـ cis-regulatory DNA منها البروموتور الأساسي (القريب) حيث يرتبط أنزيم النسخ pol II والمميز بالحروف (Bp)، والمناطق الأخرى المميزة بالحروف من A إلى G كل منها حوالى 300 bp وكل به أكثر من ١٠ مواضع حيث ترتبط العديد من البروتينات المنظمة للنسخ (المنشطة والكابتة)، كما هو مبين بشكل (٨-١).

بأستعمال أساليب الهندسة الوراثية لمعرفة وظيفة كل مقطع module من تلك المقاطع على حدة، أتضح أن المقاطع A, B, G تنشط أو تنبه نسخ هذا الجين في خلايا الـ endoderm وخصوصاً في منطقة الـ midgut من المراحل المبكرة من التشكل، وأن المقاطع F, E تكبت عمله في خلايا الـ ectoderm بينما المقاطع D, C تكبت عمله في

خلايا الـ mesoderm. وبتراكم المعلومات بهذا الخصوص اتضح أن المقطع A هو المقطع الأساسي القادر على تنبيه منطقة البروموتر الأساسي Bp، وأن باقي المقاطع تظهر تأثيرها من خلال الارتباط بتتابعات محددة داخل A ولا يمكنها الارتباط بالبروموتر الأساسي مباشرة. هذه المعلومات الشيقة دفعت Davidson وفريقه إلى تصور نظام التحكم لهذا الجين من خلال تبنى مفاهيم وأفكار الهندسة الكهربائية في تصميم الدوائر الكهربائية circuits، ولكن لم يكن هدفهم بطبيعة الحال تصميم دوائر من الأسلاك والموصلات الكهربائية، فهناك فرق جوهري حيث نتعامل هنا مع أنظمة حية. وشكل (٩-١) يمثل الدائرة الجينية endo16.

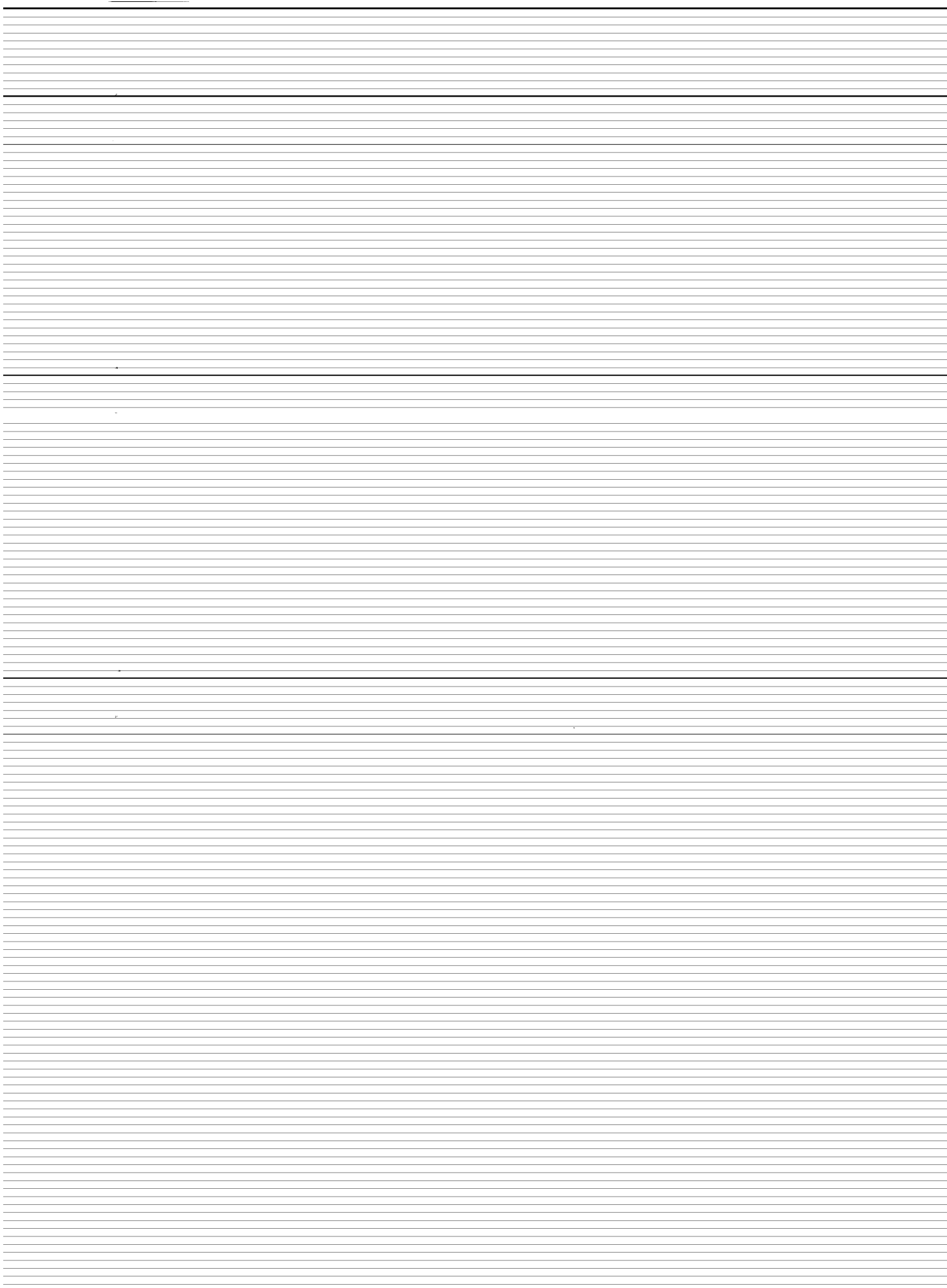


شكل (٩-١) : مناطق البروموتر البعيدة لجين endo16 مميزة لعدة مناطق حيث يمكن أن يرتبط بها البروتينات المنظمة للنسخ (الأشكال الهندسية المختلفة).



شكل (٩-١) : رسم توضيحي لدائرة الجين endo16، كما اقترحها Davidson ومعاونيه (في أقصى اليمين شكل يمثل مفتاح toggle).

وفى هذه الدائرة يمكن تصور المقاطع من A إلى G كمصادر للأشارات signals كل يخرج منه موصلات، وعند إلتقاء إشارتين مع بعضهما فتمثل على هيئة عقدة node (ممثل في الشكل بدوائر داخلها حروف يونانية) ويمكن تصور وجود ترانسيستور صغير عند كل عقدة عليه القيام بالحسابات لتأخذ القرار. فعند العقدة α يتم تلقى الإشارات من المقاطع F, E, C, D حيث ترتبط بالتتابع رقم ٥ على المقطع A، وعند العقدة β تلتقى الإشارات من التتابعين رقمي ٧ و ٨ على A لتعظيم نشاط المقطع B، وعند العقدة γ تلتقى الإشارات من التتابعات رقمي ١ و ٢ لتعظيم الإشارات من المقطع G، وعند العقدة δ تلتقى الإشارات الموجبة القادمة من المقطعين B و G وهذه العقدة متوقفة على زمن العمل، والعقدة E تلتقى الإشارات من العقد β و δ وتجمعها وتصدر إشارة موحدة من مجموعهما. أما عند العقدة ζ ، فيجب تصورها كمفتاح toggle يسمح إما بالعمل ON او القفل OFF (مثل مفتاح الكهرباء مثلا) فهو يتلقى إشارات من التتابع رقم ٦ على A (المعلومات الخاصة بالارتباط بالبروموتر الأساسي Bp لبناء النسخ) وكذلك يتلقى الإشارات القادمة من E الذي يعكس معلومات المراحل التكوينية وزمنها. هنا المفتاح زيتا ζ هو الذى يسمح بالعمل للجين، ولكن لتحديد نوع الخلايا التى يتم فيها العمل دون غيرها وجب وجود العقدة η حيث تتجمع الإشارات الكابتة من كافة المقاطع المختصة لوقف العمل فى خلايا ectoderm و خلايا mesoderm والتى يحددها موقعهما فى النسيج الجنينى. وفى حالة السماح بالعمل (خلايا الـ endoderm) تنقل الإشارات إلى θ حيث تلتقى إشارات α المساعدة على الارتباط مع Bp. ويمكننا أن تصور قدر الجهود والوقت الذى بذلته هذه المجموعة من الباحثين الأكفاء للتوصل لهذا النموذج الفريد لدائرة الجين endo16 وهو يمثل وحدة واحدة من منظومة جينومية قد يبلغ أعضائها الآلاف من الجينات، إلا أن هذا الوضع لم يثنى Davidson وفريقه سنة ٢٠٠٥ فى وضع أول تصور لشبكة جينية gene network تمثل مراحل التشكل الأول لقنفذ البحر وتشمل أكثر من حوالى ٢٠ جين مختلف، كما هو مبين بشكل (١-١٠). ومع أن هذه الشبكة مبدئية ويشوبها الكثير من القصور إلا أن المستقبل يبشر بنتائج طيبة فى هذا المجال.



٢. الجينومكس التركيبي

Structural Genomics

إعداد: أحمد المتينى

لكي نقرب إلى ذهن مدى الجهد والصعاب التي قابلها المختصين لإنجاز الدراسات والمشاريع لسلسلة أو فك تتابعات جينومات بعض الكائنات والتي انتهت، خصوصا في مراحلها الأولى، يمكن أن نجرى المقارنة التالية، فلو أخذنا مثلا الفيروس (البكتريوفاج لامدا) فجينومه حوالى ٥٠٠٠٠ bp فإذا افترضنا أن الصفحة المكتوبة بخط صغير جدا تحتوى على تتابعات قدرها ٢٥٠٠٠ bp ، فإن جينومه بالتالى سيملا صفحتين. إذا انتقلنا إلى البكتيريا (*E. coli*) فإن تتابعاتها ستشغل كتيباً صغيراً من ٢٠٠ صفحة، فإذا انتقلنا إلى كائن دقيق حقيقى النواة مثل الخميرة فإن تتابعات جينومه ستملأ مجلدا صغيراً من حوالى ٥٠٠ صفحة. أما الدودة النعمانية (*C. elegans*) والنبات الصغير (*A. thaliana*) فهما جينومين متساوين تقريبا فى الحجم، فجينوم كل منهما سيشغل ثلاث مجلدات كبيرة، أما إذا وصلنا للإنسان فإن جينومه سيشغل حوالى ٨٠ مجلدا كبيرا. والجدول التالى يلخص قدر أزواج النيوكلووتيدات فى مجاميع الكائنات الرئيسية.

مجموعة الكائنات	حجم الجينوم (bp)
الفيروسات	300 – 350,000
غير حقيقية النواة	250,000 – 15,000,000
حقيقية النواة – وحيدة الخلية	12,000,000 – 50,000,000,000
حقيقية النواة – عديدة الخلايا	20,000,000 – 600,000,000,000

وبزيادة المعلومات لعلم الجينومكس وتشعبها اتفق العلماء على تقسيمه إلى عدة

مجالات وهي :

(١) الجينومكس التركيبي (Structural genomics) الذى يختص بدراسة التركيب

الجزيئي (المادي) لجينوم ما ،

(٢) الجينومكس الوظيفي (Functional genomics) الذى يختص بالتعرف على النشاط

أو الفعل لجينوم ما،

(٣) الجينومكس المقارن (Comparative genomics) ويختص بمقارنة جينومات الأنواع

المختلفة بعضها ببعض.

١.٢. مدخل Introduction

عادة يحدد حجم جينوم الكائن مدى الجهد والصعوبة المتوقعة لدراسته، لكن يجب أن نوضح بأن حجم الجينوم لكائن ما لا يعكس بالضرورة مدى تطور هذا الكائن فيما يعرف باسم معضلة حجم الجينوم (C-paradox) فبعض البرمائيات وبعض النباتات تتميز بجينوم (C-value) أكبر من ذلك فى الإنسان. ولدراسة تركيب جينوم ما يتم أولا عمل خرائط جينية له ومنها تعمل خرائط جزيئية ثم يتم كونه كل قطعة متداخلة (contigs) والتي فى مجموعها تغطى جزء من الجينوم ويتبع ذلك فك تتابعات كل منها. أي يجب هدم الجينوم إلى وحدات أصغر فأصغر حتى نصل إلى قطع تسمح للأساليب العملية بفك تتابعاتها فيما يسمى [Top-Down approach for molecular mapping] وشكل (١.٢) يوضح تلك الخطوات كما اقترحها Klug & Cummings سنة ١٩٩٩. وإذا ما وصلنا إلى هذه النقطة، وجب علينا أن نعكس العمل وأن نعيد البناء مستعملين تلك القطع الصغيرة لنصل إلى الأكبر فالأكبر حتى الجينوم فيما يعرف باسم [Bottom-Up approach].



شكل (١-٢) : رسم يوضح خطوات دراسة جينوم ما.

٢.٢. الخرائط الجينومية Genomic maps

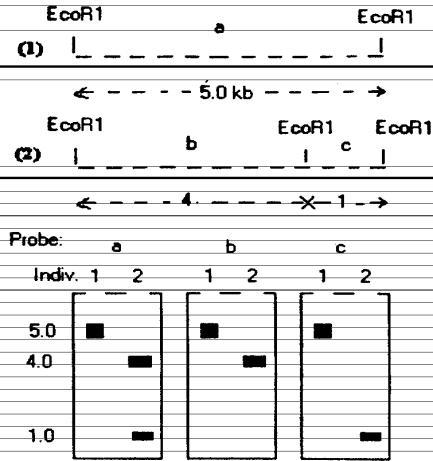
عرف علماء الوراثة الخرائط الجينية منذ بداية العشرينات من القرن الماضي، وكانت تعتمد على نسب العبور بين المواقع الجينية (loci) المختلفة التي تتعدد عليها الاختلافات الظاهرية (phenotypic variants). هذه الاختلافات (الطفرات) عند جمعها بكتائن ما، فإنها تكون عادة مقللة للحيوية أو حتى مميتة، كما أن هذه الطريقة تحتاج لعدد كبير من التهجينات - أو دراسة عدد كبير من سجلات العائلات الإنسانية. هذه الأسباب قللت من إمكانية استخدامها على نطاق واسع مع العلم بأن فترة هذه الخرائط لا تسمح إلا بتحديد مواقع على الخريطة تبعد عن بعضها البعض بأكثر من 1 Mb، ولا يتوفر الآن سوى عدد قليل من البرامج التي تساعد في عمل هذا النوع من الخرائط مثل برنامج MAPMAKER المتوفر بواسطة معهد Whitehead Inst., MIT بالولايات المتحدة الأمريكية. وظل الحال على ما هو عليه دون تقدم ملحوظ إلى أن اكتشفت للسومات الجزيئية (molecular markers)، حيث تغير الحال وامتدنا بعدد كبير من التسلسلات الجينية على مستوى تنافس DNA حيث تتميز بتعدد الأشكال الظاهرية (polymorphic) كما أنها متعادلة (neutral) التأثير على حيوية الفرد وتنتقل تبعاً

للقواعد المندلية. استعمال تلك السمومات دفع بدراسة الخرائط الجينومية دفعة كبيرة إلى الأمام مكنتنا من زيادة قوة الإيضاح إلى حوالي 100 kb ولكن مثل تلك الأحجام من الجينوم مازالت كبيرة وغير مناسبة لفك تتابعاتها، لذلك وجب اللجوء إلى أساليب إضافية للمساعدة في هذا المجال.

١.٢.٢. عمل الخرائط ببصمة إنزيمات القصر. Restriction fingerprint mapping.

عادة تقطع أجزاء الـ DNA إلى شظايا صغيرة بواسطة عدد من إنزيمات القصر لتكون مناسبة لربطها بأدوات النقل (cloning vectors) مثل اليكتروفاج لامدا الذي يحمل حوالي (20 kb) أو الكوزميدات تحمل حوالي (40 kb) أو الـ YAC الذي يحمل حوالي (200 kb)، ومع استخدام المجسات الجينية (probes) للتعرف على تلك الشظايا. فعلى سبيل المثال لو استعملنا أنزيم القصر EcoR1 لقطع قطعة كبيرة من الـ DNA فإن أحجام تلك القطع سيتوقف على أماكن تعارف هذا الأنزيم كما هو مبين بشكل (٢-٢)، ففي فرد ما (رغم ١) ستنتج قطعة قدرها 5 kb لوجود موقعين للقطع بينما وفي فرد آخر (رقم ٢) ستنتج قطعتين قدر الأولى 4 kb والثانية 1 kb نتيجة لوجود ثلاث مواقع لقطع الأنزيم في ذات الموقع (تعدد أشكال ظاهرية). لعزل تلك القطع والتعرف عليها يمكن استخدام المجسات الجينية (probes) ولتكن في هذه الحالة المجسات a و b و c التي تمثل الشظايا المتوقعة الحصول عليها بعد تعليمها إما بالنظائر المشعة أو بمركب فلوروسنتي. وعند الفصل الكهربائي للـ DNA على الأجاروز وبالتجهين مع المجسات تبعا لطريقة Southern blotting فستظهر الحزم fragments المبينة بشكل (٢-٢) حسب الفرد والمجس المستعمل.

هذه الشظايا العديدة التي تم كوليبتها وفك تتابعاتها تقطى جزء من جينوم ما ولكنها متداخلة مع بعضها البعض (overlapping) لذلك تعرف بشظايا الـ DNA المتداخلة أو ما يسمى بـ contigs. وإذا نجحنا في تحديد أماكن بعض منها على الجينوم فإنها تصبح نوعا خاصا من المجسات الجينية تعرف بالتتابعات المحددة لموقع بعينه (sequence tagged sites) وتعرف اختصارا باسم STS.



شكل (٢-٢) : رسم يوضح خطوات قطع الـ DNA بانزيم القصر EcoRI مع تمييزه

باستخدام المجسات probes لجين ما .

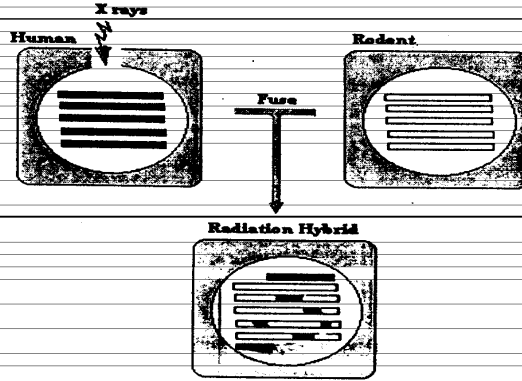
٢.٢.٢. طريقة دمج الخلايا بعد الإشعاع Radiation Hybrid (RH)

استغل العلماء النجاح في التهجين أو الدمج ما بين الخلايا الجسمية (somatic cell hybridization) واستعملوها لتحديد أماكن الجينات على الكروموسومات. وقد استعملت خلايا لإنسان بعد إشعاعها بجرعة من أشعة X قدرها ٢٠٠٠ راد لتكسير الكروموسومات إلى شظايا صغيرة، ثم تدمج (fuse) تلك الخلايا المشعة مع خلايا الفأر، كما هو مبين بشكل (٢-٢). وبفحص الخلايا المهجنة (الدمجة) يلاحظ ارتباط شظايا إنسانية بكروموسومات الفأر بالإضافة لقطع غير مندمجة، وبحساب نسب تكرار الاندماجات المفردة والمزدوجة والتي تمثل الشظايا المتجاورة يمكن تقدير المسافات بينها، حيث وجد أن وحدة centi Rays أو (cR3000) تساوي 0.1 cM لهذا النظام. وقد اقترح Griffith ومساعدوه سنة ١٩٩٩ أن استعمال من ١٠٠ إلى ٢٠٠ خلية مدمجة من هذا النوع ستكون كافية لعمل خرائط قدرة إيضاها عشرة أضعاف الخرائط الوراثية التقليدية.

٢.٢.٢. طريقة استعمال المجسات على المستوى الخلوي.

Fluorescence *In situ* Hybridization (FISH).

أن استخدام طريقة التهجين على المستوى الخلوي بالمجسات الجينية المعلقة بمرسبات فلوروسنتية، وتعرف اختصاراً باسم FISH، يمكن استخدامها أيضاً لتحديد بعض المواقع على كروموسومات الوضع المتوسط (metaphase)، كما هو موضح بشكل (٤-٢).



شكل (٢-٢) ، رسم توضيحي لطريقة RH في الإنسان.



شكل (٤-٢) ، صورة لكروموسومات الوضع المتوسط للإنسان توضح ارتباط مجس محدد على

موضع بكروموسوم ١٧ باستخدام طريقة FISH.

والآن يمكن استخدام أكثر من مجس واحد للتحضير الواحد وذلك عند استخدام مجسات معلمة بمركبات فلوروسنتية ذات ألوان مختلفة بما يسمى chromosome painting. وقد اتضح أن استعمال كروموسومات الوضع المتوسط يؤدي إلى خرائط ضعيفة الوضوح ويرجع ذلك إلى مستوى التكاثف (condensation) والالتفاف الذي يتعرض له الـ DNA خلال هذا الدور من الانقسام، مما لا يتيح تحديد للمواقع التي يفصلها عن بعضها البعض مسافات تقل عن عدة ملايين من القواعد. ولزيادة كفاءة هذه الطريقة اقترح استخدام خلايا الدور البيني (intrapase)، كما طورت الطريقة بعد ذلك عن طريق فرد خيوط الـ DNA كيميائياً في التحضيرات السيتولوجية قبل تثبيتها وصيغها فيما يسمى بطريقة (fiber-FISH) طبقاً للعالم Hliskanen و مساعدوه سنة ١٩٩٥ وبذلك أمكنهم تحديد مواقع لا يبعدها عن بعض سوى 100 - 500 kb.

٢.٢. المكتبات الجينومية. Genomic libraries

اتفق العلماء على أن مصطلح المكتبات الجينومية، يقصد به مجموعة الكلونات (clones) التي تمثل جزءاً اعتباطياً من الجينوم. وهذا العمل يحتاج إلى مجهود مضني مستخدمين العديد من التقنيات الحديثة، وفيما يلي الخطوط العريضة التي يجب إتباعها لعمل مكتبة لجينوم خاص بكانن ما.

١.٢.٢. تحديد عدد الكلونات المطلوبة.

Determining the number of clones needed.

أن تحديد عدد الشظايا (fragments) من الـ DNA التي يجب كlonتها يعتبر من المهام التمهيدية الهامة التي يجب الاهتمام بها قبل البدء في عمل المكتبة. فأن استخدامنا عدد قليل من الكلونات فإنه بالتبعية لن يغطي القطعة الجينومية بالكامل، وإن استخدمنا عدد كبيراً من الكلونات فإنها ستغطي مناطق من القطعة الجينومية عدة مرات (تكرار) مما يعتبر إهداراً للجهد والمال. والمعادلة التالية توفر لنا حساب عدد الكلونات المطلوبة.

$$N = \frac{\ln(1-P)}{\ln(1-f)}$$

N = عدد الكلونات المطلوبة للحصول على كلون واحد لكل جين على هذه القطعة

الجينومية عند مستوى احتمال P قدره ٩٩٪،

f = النسبة لم توسط حجم الشظية بالكلون إلى متوسط حجم الجينوم الكلي.

فعلى سبيل المثال: أنا كلن متوسط الشظية التي يحملها الفاج لامدا حوالي 0.02 Mb، ومتوسط حجم الجينوم الكلي لنبات الأرابيدوبسيس حوالي 70 Mb فإن:

$$f = 2.86 \times 10^{-4} \quad \& \quad N = 1.61 \times 10^4.$$

ويمكن تبسيط تلك القيم، بالقول بأن تقسيم جينوم الأرابيدوبسيس إلى شظايا قدر كل منها 20kb بحيث يمثل كل جين بكلونة واحدة سنحتاج إلى ٢٤٩٧ كلونة في حالتنا هذه. أي لتحديد كل جين على هذا الجينوم عند مستوى احتمال قدره ٩٩٪ فعلى فحص عدد 1.61×10^4 كلونة

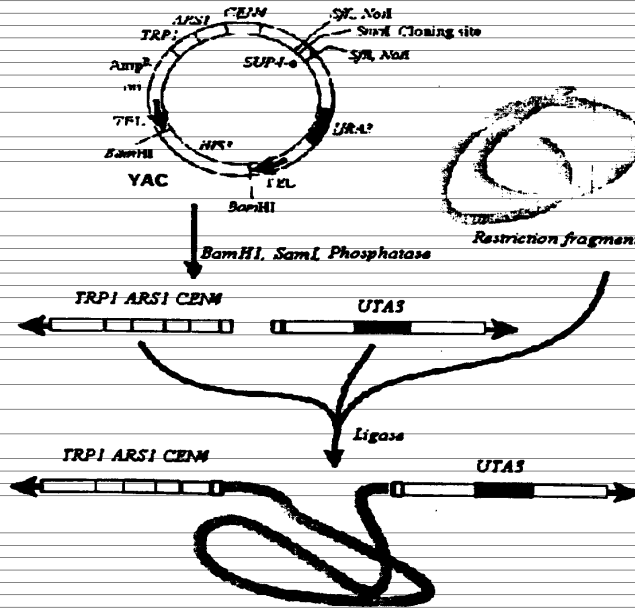
٢.٢.٢. أدوات الكولنة. Cloning Vehicles

تقدم لنا تقنيات الوراثة الجزيئية العديد من الأدوات والأساليب التي تساعدنا في كولنة الشظايا الجينومية داخل أدوات للنقل لعمل المكتبات الجينومية، ولزيد من التفاصيل في هذا الصدد يمكن الرجوع إلى (أبويوسف و المتينى - ٢٠٠٢). وبطبيعة الحال فكلما زاد حجم القطع المكونة كلما قلت الحاجة لأعداد أكبر من الكلونات لتغطية الجينوم تحت الدراسة. ويعتبر كروموسوم الخميرة الصناعي (yeast artificial chromosome) والمعروف اختصاراً باسم YAC من أنسب الأدوات لهذا الغرض لقدرته على حمل شظية كبيرة من ١٠٠ - ١٠٠٠ قاعدة (في المتوسط ٢٠٠ قاعدة)، كذلك فإن YAC يمكن أن يتضاعف كبلازميد في البكتيريا *E. coli* و كروموسوم في الخميرة. ويتميز بتتابعات خاصة مثل منطقة السترومير [CNE] و منطقة التضاعف الناتى [ARS] و منطقة التلومير التي يمكن فتحها بأنزيم BamH1 [TEL] و العديد من المواقع التي تساعد في إنتخاب الكلونات مثل [His3 , URA3 , TRP] بالإضافة

لمواقع التعايش في البكتيريا [Amp, ori]. شكل (٥-٢) يمثل رسماً توضيحياً لطريقة ربط قطعة جينومية مع الـ YAC.

٢.٢.٢ ترتيب المكتبات. Ordering of Libraries

المكتبة الجينومية التي حصلنا عليها حتى الآن عبارة عن كلونات موزعة اعتباطياً (غير مرتبة) لجينوم ما، لذلك يجب علينا الآن العمل على ترتيب تلك الكلونات بالبحث أساساً عن مناطق التداخل بينها أو بما يسمى بالـ contigs - أي الشظايا المكونة بها مناطق متماثلة مع تمثيل المناطق غير المتداخلة بينها مرة واحدة فقط. وبترتيب المكتبة نكون قد خطونا خطوة كبيرة إلى الأمام، فمنها يمكن أن نحدد أماكن الجينات المرغوبة وعدد نسخ كل منها بل يمكننا أن نعرف الكثير عن التعقيد التركيبي للجينوم ككل.

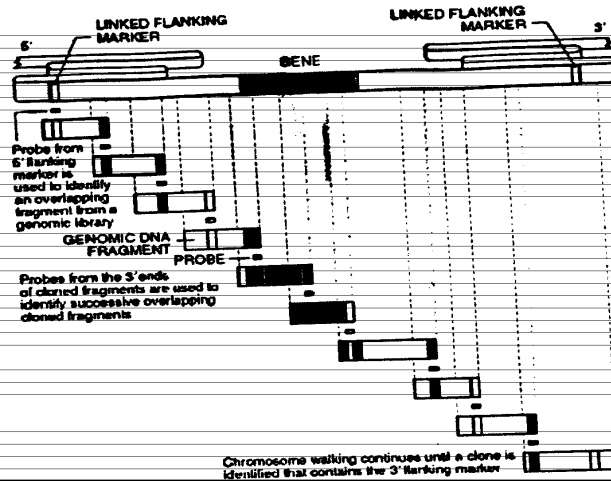


شكل (٥-٢) : رسم توضيحي لأسلوب ربط شظية جينومية مع الـ YAC كأداة للكلونة.

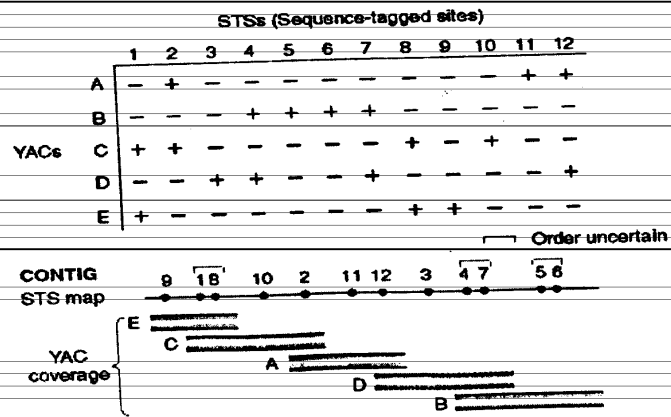
وأبسط الطرق لتحقيق ذلك يتم بواسطة البدء بكلونة ما من المكتبة تحت الدراسة، ثم البحث عن ثمانية بها تتابعات متماثلة مع بعض من القطعة الأولى، ثم البحث عن قطعة ثالثة لها تتابعات متماثلة مع بعض من تلك بالقطعة الثانية ... وهكذا. وأول الطرق التي استغلت لتنفيذ ذلك وتعرف بأسم المشي الكروموسومي (Chromosome walking) أو القفز الكروموسومي (Chromosome jumping) وخطوات المشي الكروموسومي موضحة بشكل (٦-٢). يتبع الآن طرق عديدة أخرى خلافاً للمشي الكروموسومي لترتيب المكتبات تعتمد على التهجين (hybridization) بين الجزيئات وعمل البصمات (fingerprints) للـ DNA عن طريق هدمه إلى قطع صغيرة، ولابد من استعمال برامج الكمبيوتر لأجراء آلاف المقارنات ما بين الكلونات لتحديد مناطق التداخل وترتيب الـ contigs. ومن الطرق شائعة الاستخدام طريقة بطاقات التتابعات المحددة الموقع على الجينوم sequence tagged sites أو ما يعرف اختصاراً STS. وهي تتابعات صغيرة جداً يمكن تحديد مكانها على الجينوم، فمثلاً إذا كانت الشظية (١) عرف عليه الطاقتين ١ و ٢ والشظية (ب) عرف عليها البطاقتين ٢ و ٣، بالتالي فإن الشظيتين أ و ب متداخلتين عند البطاقة ٢ وهكذا. وشكل (٧-٢) يوضح تلك الخطوات تبعا لـ Griffith ومساعدوه سنة ١٩٩٩.

٤.٢ فك (سلسلة) التتابع Sequencing.

تقنية فك (أو معرفة) التتابع هي تلك الأساليب العملية المتبعة لتحديد تتابعات النيوكليوتيدات الأربع على طول شظية صغيرة من الـ DNA. والأساس الذي يبنى عليه ذلك هي الطريقة الكيماوية الأنزيمية التي اقترحها Sanger منذ أكثر من ثلاثين عاما والتي ما تزال تستعمل حتى الآن وإن أدخلت عليها تعديلات لزيادة الكفاءة والاعتماد على التنفيذ الأوتوماتيكي. وتستعمل فيها النيوكليوتيدات الأربع معلمة بالنظائر المشعة أو هرواند الفلوروسيننتية مع النيوكليوتيدات المعدلة (dideoxynucleotides-ddN) التي توقف عمل أنزيم البوليميراز عند ارتباطها وبالتالي تحدد النيوكليوتيدة الأخيرة وهكذا عن طريق فصل الشظايا تبعا للحجم بالتفريد الكهربائي، كما هو مبين بالشكل شكل(٨-٢). كما تتوفر الآن أجهزة إلكترونية مبرمجة يمكنها فك وتحليل التتابعات مباشرة وتقديم النتائج كملاقات بيانية، كما هو مبين بشكل (٩-٢).



شكل (٦-٢) : خطوات عملية المشي الكروموسومي.



شكل (٧-٢) : رسم يوضح طريقة استعمال الـ STS لتحديد ترتيب الـ contigs.

عادة لسلسلة تتابعات شظايا جينومية ما نلجأ إلى استراتيجيتين أساسيتين، الأولى تعتمد على فك تتابع القطع المكونة المرتبة (ordered clone sequencing) ، والثانية هي فك تتابع الجينوم ككل (whole genome shotgun sequencing) والطريقة الثانية متبعة على نطاق واسع مع الجينومات الصغيرة للكائنات غير حقيقية النواة.

١.٤.٢ Ordered Clone Sequencing

ويمكن تلخيص هذه الاستراتيجية في الخطوات التالية:

١. عمل مكتبة جينومية مستغلين أداة نقل تتحمل شظايا كبيرة نسبياً.
٢. عمل خريطة جزيئية مرتبة لهذه المكتبة.
٣. عمل ما يسمى بخط السير الأدنى (minimum tilling path) ويقصد به العدد الأدنى الممكن من الـ contigs التي تغطي الجينوم بأكمله.
٤. إعادة كونة هذه الـ contigs إلى شظايا أصغر حجماً (2 kb) مستعملين الأداة المناسبة.
٥. فك تتابع الشظايا الصغيرة باستعمال مناطق الارتباط مع الأداة كبادئات (primers) ويفك التتابع في كلا الاتجاهين بدءاً من تلك المواقع.
٦. إعادة تجميع الشظايا (معروفة التتابعات) في contigs مرة أخرى.

٢.٤.٢ Whole Genome Shotgun Sequencing

ويمكن تلخيص هذه الاستراتيجية في الخطوات التالية:

- عمل مجموعة عشوائية من الشظايا الصغيرة لجينوم ما عن طريق الهدم الميكانيكي (sonication) أو الهدم الأنزيمي ولكن تفضل الطريقة الأولى لتجنب التداخل الأنزيمي.
- التعرف على تتابعات أكبر عدد ممكن من تلك الشظايا الجينومية، يمكن الاكتفاء بفك تتابع أطراف الشظايا فقط فيما يسمى (long insert library).

- تجمع تلك الشظايا (معروفة بالتتابع) في contigs متداخلة، ولكن عادة ستبقى فجوات (gaps) فيما بينها.
- ترتب contigs بالنسبة لبعضها البعض.
- تطلق الفجوات المتبقية تبعا لطريقة primer walking.

٥.٢. بطاقات التتابعات المميزة للجينات الفعالة.

Expressed Sequence Tags [ESTs]

بطاقات التتابعات الفعالة والتي تعرف اختصاراً بـ ESTs هي تتابعات من الـ DNA المبنية من جزيئات الـ DNA المكمل (cDNA) لرسالات الجينات التي عبرت عن فعلها بالنسخ في نسيج ما تحت ظروف محددة. وكان العالم البيولوجي الشهير Venter قد اقترح في أواخر الثمانينات من القرن الماضي على المعهد القومي الأمريكي للصحة (NIH) مشروعاً بحثياً يدعو إلى التركيز على التعرف على الجينات النشطة فقط في الجينوم الأنساني عن طريق التعرف على تتابعات مميزة للجينات النشطة فقط والمستخلصة من مرسالات (mRNA) الجينات التي نسخت بالأنسجة المختلفة، فيما أطلق عليه اسم ESTs، وبرر ذلك بأن الجينات الفعالة لا تمثل أكثر من ٢٪ من جملة الجينوم الأنساني مما يجعل البحث عن جين ما مهمة صعبة للغاية، وأن تحديد موقع الجين في حد ذاته لا يعكس طبيعة نشاط هذا الجين، كذلك فهو توفير للمال والجهد. ولكن رفضت الهيئة الأمريكية تمويل المشروع في حينه مما دفع Venter إلى اللجوء إلى مجموعة من رجال الأعمال واهتمهم بتمويل إنشاء هيئة خاصة هدفها الأساسي التعرف على تلك البطاقات بالجينوم الأنساني ومحاولة استغلالها في المجال الطبي، وتبعاً لذلك أنشأ معهد أبحاث الجينوم (The Institute for Genome Research) والمعروف اختصاراً TIGR <http://www.tigr.org> والذي أصبح بعد ذلك من أشهر الهيئات في هذا المجال. ومنذ ذلك الوقت أنتجت الـ ESTs على نطاق كبير، كانت تمثل الجينوم الأنساني في أول الأمر ثم ما لبس أن تبعه العديد من الكائنات الأخرى، فحتى أكتوبر من سنة ١٩٩٩ تم تسجيل في بنك الجينات (GenBank - <http://www.ncbi.nlm.gov>) حوالى ثلاثة ملايين بطاقة لعدد من الكائنات. وفي سنة ١٩٩٥ قام العالم Adams مع أكثر من تسعين باحثاً في نشر قائمة بهذه البطاقات في الإنسان في مجلة Nature زائفة الصيت مثلت أكثر من ٨٠

الف جين إنساني ٨٠٪ منها لم تكن معروفة من قبل، والجدول التالي يلخص هذه البطاقات لكل عضو درس حتى تاريخه.

مصدر الـ DNA	عدد البطاقات	مصدر الـ DNA	عدد البطاقات
الجلد	٢٠٤٢	خلايا الدم	٢٥٠٥
العظم	٥٧٢٨	الغدة الدرقية	٢٢٧٨
الغدة جارة الدرقية	١٩٧	الأغشية المصلية	٥٧٣٦
الكبد	٢٧٨٠٧	المرارة	٢٧٥٤
القولون	٤٨٣٢	العضلات الملساء	٢٩٧
الأمعاء الدقيقة	١٠٠٩	البروستاتا	٧٩٧١
الخصية	٧١١٧	المبيض	٣٨٤٨
الرحم	٦٣٩٢	المشيمة	١٢١٤٨
المخ	٦٧٦٧٩	العين	١٩٣٢
الغدة اللعابية	٨٨	المريء	١٩٤
الأنسجة الدهنية	٢٤١٢	القلب	٩٤٠٠
الغشاء البروتوني	٨٨٤	الطحال	٧٩٢٤
الکظر	٢٤٢٧	البنكرياس	٥٥٣٤
الکلي	٢٢١٢	البربخ	١٧١٦
العضلات	٤٦٩٢	الأغشية الزلالية	٢٨٨٩
الجنين	١٩٢٩١	الغدة التيموسية	٢٤١٢

لتوضيح أهمية هذه البطاقات ESTs في تحديد الجينات، دعنا ندرس مثال على ذلك، فبفرض أننا حددنا موقع جين ما بين الموقعين ٨٧ و ٩٣ وحدة عبورية سنتمورجان (cM) على الكروموسوم رقم ١٦ في الإنسان، فيمكننا الرجوع إلى الجداول الخاصة حيث يمكن اختيار بعض البطاقات ESTs التي تتواجد ما بين هذين الموقعين وفي أي الأعضاء هي نشطة، فوجد أن البطاقة رقم A006F26 ترتبط بهذا الموقع وهي

في الأصل تظهر تماثلاً كبيراً مع الجين الخاص بالحالة المرضية المعروفة باسم

. Huntington's Chorea

٣- الجينومكس الوظيفي

Functional Genomics

إعداد : ياسر مبروك و أحمد المتيني

مع بدأ العمل في مشاريع فكّ وتحليل الجينومات للعديد من الكائنات الأولية والراقية، ظهرت الحاجة إلى ابتكار وتوفير الأساليب العملية وكذلك الأسس النظرية لمعرفة وتحديد وظائف وطبيعة عمل الجينوم ككل، وعليه ظهر هذا الفرع الجديد من العلوم البيولوجية الجزيئية والذي يعرف باسم الجينومكس الوظيفي Functional genomics. ومما سبق يتضح لنا أن هذا الفرع يهتم بدراسة وتحديد أماكن النشاط الجيني في الجينوم، ومن خلال هذه الدراسة يتوقع منا محاولة الإجابة على عدة تساؤلات :

(١) أين ومتى يتم التعبير عن جين ما ؟

(٢) ماهي طبيعة ناتج الفعل الجيني ؟

(٣) ما هو مدى التعاون بين الجينات لإظهار فعلها ؟

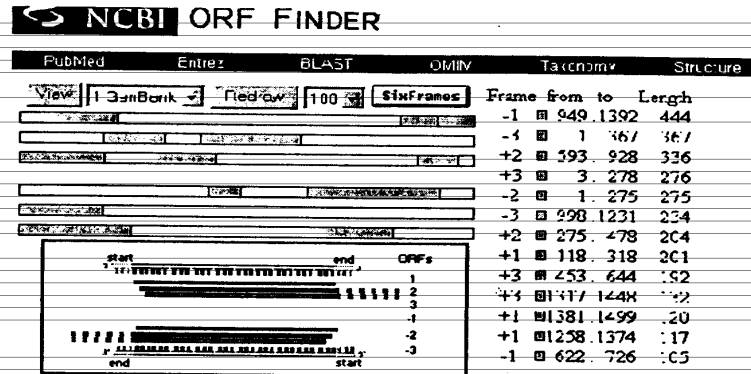
في محاولة أجابنا على السؤال الأول، لابد ان نتناول بالتدقيق مواضيع محددة مثل ما هي مواصفات المواقع التي تنسخ؟ وما دور عوامل التنظيمية وتداخلاتها على المناطق غير المنسوخة؟ هذا المجال من الدراسات يمكن إجماله فيما يسمى بال- Transcriptomics = بمعنى آخر هو دراسة الرسائل mRNAs. أما التساؤل الثاني فسيختص بدراسة عمليات الترجمة وتكوين الناتج الجيني، ألا وهو البروتينات. هذا المجال من الدراسات يمكن إجماله بما يسمى بال- Proteomics – وسوف نتناوله بشيء من التفصيل فيما بعد. أما تساؤلنا الثالث فيهتم بدراسة التداخل والتعاون بين الجينات لأظهار الشكل الظاهري فيما يسمى ب- Phenomics – وقد نتناول هذا الموضوع بالدراسة من خلال مواضيع الوراثة التكوينية Developmental genetics.

١.٢. التنبؤ بمواقع الجينات. Prededction of Gene Locations.

إن الإطارات المفتوحة للترجمة (ORFs – open-reading frames) هو تعبير مرادف لمصطلح الجينات النشطة التي تنسخ وترجم. وعادة يكون عدد الـ ORFs في جينوما ما محدودة للغاية مقارنة بباقي المادة الوراثية بالجينوم، ومع ازدياد عدد الجينومات التي مسحت وفحصت كان من الضروري محاولة البحث أو التنبؤ بهذه المواقع فيما يسمى ORFs prediction والتي لم تكن معروفة من قبل. وعادة لعمل هذا البحث نلجأ لبرامج الكمبيوتر المختصة اعتماداً على خواص الـ ORFs أو ما يطلق عليه مسمى الـ signals، فعادة تبدأ بكودون "بدء الترجمة" وهو غالباً AUG ولا بد أن ينتهي بأي واحدة من الكودونات الثلاثة الخاصة بالإنهاء، ويجب أن يكون طویل نسبياً (على الأقل ١٠٠ كودون) حتى يترجم إلى بروتين ذو معنى. وعلى ذلك فطول الـ ORF الافتراضي سيساوى $2/64$ حيث ٦٤ هو عدد الكودونات الكلى و٢ هو عدد كودونات الإنهاء. وفي أغلب الأحيان تستعمل برامج الكمبيوتر المبنية على النماذج الإحصائية خصوصاً نماذج ماركوف (سنناولها ببعض التفصيل بالباب القادم) للتنبؤ وتحديد هذه الإطارات من ضمن الإطارات الستة المحتملة لكل ORF كما هو موضح بشكل (٢-١). والاعتماد على النماذج الإحصائية في كثير من الأحيان يؤدي إلى استبعاد معطيات هامة : مثل الأكسونات الصغيرة والتي يقل فيها عدد الأحماض الأمينية عن ١٠٠، كذلك تشخيص إطارات يزيد طولها عن ١٠٠ ولكنها تدرج تحت نوعية الجينات التي تنسخ ولا تترجم..... وهكذا العديد من الحالات المتداخلة، لذلك وجب تأكيد تقديرات الكمبيوتر وحساباته بأجراء التجارب العملية المبنية على تقنيات البيولوجيا الجزيئية.

مما سبق، فإن هذا المفهوم يكون الاتجاه العام الذي سنتبعه في تناولنا لموضوعات الجينومكس والمعلوماتية الحيوية. فلدراسة الفعل الجيني تحت مسمى الجينومكس الوظيفي يمكن أن نسلک أحد مسارين أو كليهما معا للوصول للهدف المرجو، الأول من خلال تطبيقات علوم الكمبيوتر واعتماداً على قواعد البيانات المتخصصة من خلال البحث والمقارنة والاستدلال والتنبؤ فيما يسمى بالبحث عن القرائن homology، وهي ما سنتناوله في الأبواب التالية – أو ما يعرف بالمسار المعلوماتي. أما المسار الثاني – ألا وهو

المسار التجريبي تبعا لتقنيات البيولوجيا الجزيئية، فهو يشمل عدة اتجاهات يمكن إجمالها في بعض الأساليب التالية.

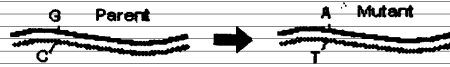


شكل (٣-١) : رسم توضيحي لأحد الجينات توضح الأطارات الستة المحتملة له بناء على أن الكودونات تتكون من ثلاث نيوكليوتيدات. بالإضافة لصفحة من برنامج ORF Finder المتوفر من قاعدة NCBI يوضح نتائج البحث عن هذه الأطارات لجين ما.

٢.٢. أساليب الإطاحة بفعل الجين Gene Knockout.

هذا الاتجاه يعتبر من أكثر الاتجاهات شيوعا لدراسة الفعل الجينومي، فعن طريق إسكات عمل جين ما يمكن تتبع التفاعلات التي ستظهر على الشكل الظاهري حيث سمى هذا الاتجاه بالوراثة العكسية reverse genetics لكونها عكس ما جرى عليه الحال في الدراسات الوراثية التقليدية. ويرجع الفضل إلى العالم Michael Smith (نال جائزة نوبل للعلوم سنة ١٩٩٢) في أواخر السبعينات من القرن الماضي لتطوير طريقة لاستحداث الطفرات داخل النظام الحي *in vivo* mutagenesis والمعروف أيضا باسم الطفرور موجه لموقع site-directed mutagenesis لاستحداث طفرة موضعية عند موقع محدد من تتابع ما من الـ DNA قد يمثل ORF. وقد طورت تقنيات إستحداث الطفرات هذه فيما بعد وأصبحت الطرق الحديثة أكثر كفاءة حتى وصلت احتمالات النجاح الآن لأكثر من ٧٥% في الحصول على الطفرة المرغوبة، فلو كان عندنا قطعة من الـ DNA تمثل ORF

ما (قدعنا نعرفها بالجزئ الأبوي parent molecule) وإذا ما كنا نرغب في إدخال طفرة (استبدال نيوكليوتيدات) على موقع محدد فيمكن تمثيلها كالتالي:



وهذا الاستبدال قد يكون عند الأطراف أو عند منتصف الجزئ.

١.٢.٣. إستبدالات substitutions عند أطراف الجزئ.

بأستعمال تقنيات تفاعلات البلمرة المتسلسلة polymerase chain reactions

المعروف اختصاراً بـ PCR يمكن أحداث مثل هذه التغيرات. فمثلاً لو هناك قطعة من الـ DNA تمثل تتابع ما كما هو مبين بالتالي:

5'TCTATGGACCAGTACGATACCGTA.....CGACCTACGTAGACTAGACGGATAGAG 3'
3' AGATACCTGGTCATGCTATGGTCAT..... GCTGGATGCATCTGATCTGCCTATCTC 5'

وهذا هو التتابع الأصلي (الأبوي) parental molecule، فلتغير أطرافها يمكن أستعمال زوج من البادئات primers لتعظيمها amplification بواسطة الـ PCR، ويجب أن يكون تصميمها كالتالي:

البادئ من الشمال = 5' TCTATGGACCAGTACGAT 3'

والبادئ من اليمين = 5' CTCTATCCGTCTAGTCTA 3'

فبأستعمال هذين البادئين يمكن تعظيم تلك القطعة ملايين المرات، ولكن لو

أريد أن تتضمن تلك القطعة تتابعات لبعض إنزيمات القصر مثل EcoRI (GAATTC) من الطرف الشمال ومثل BamHI (GGATCC) من الطرف اليمين فيمكن إضافة تلك التتابعات لكل من البادئين السابقين ليصبح البادئ الأيسر

5' GCGAATTCTCTATGGACCAGTACGAT 3'

والبادئ الأيمن 3' GCGGATCCCTCTATCCGTCTAGTCTA 5'

ويلاحظ إضافة تتابع GC من كلا الطرفين وذلك لتسهيل عمل القطع
للأنزيمين فيما بعد، وبإجراء الـ PCR مرة أخرى باستعمال هذين البادئين فإن القطعة
المعظمة amplified fragment ستصبح كالتالي:

GCGAATTCTCTATGGACCAGTA...
CGATACCAGTACGACCTACGTAGACTAGACGGATAGAGGGATCCGC
CGCTTAAGAGATACCTGGTCATGCTATGGTCAT.....
GCTGGATGCATCTGATCTGCCTATCTCCCTAGGCG

ويقطع هذه الجزيئات بأنزيمي القصر EcoRI و BamHI سنحصل على القطع
التالية:

AATTCTCTATGGACCAGTACGATACCAGTA...
CGACCTACGTAGACTAGACGGATAGAGG
GAGATACCTGGTCATGCTATGGTCAT...
GCTGGATGCATCTGATCTGCCTATCTCCCTAG

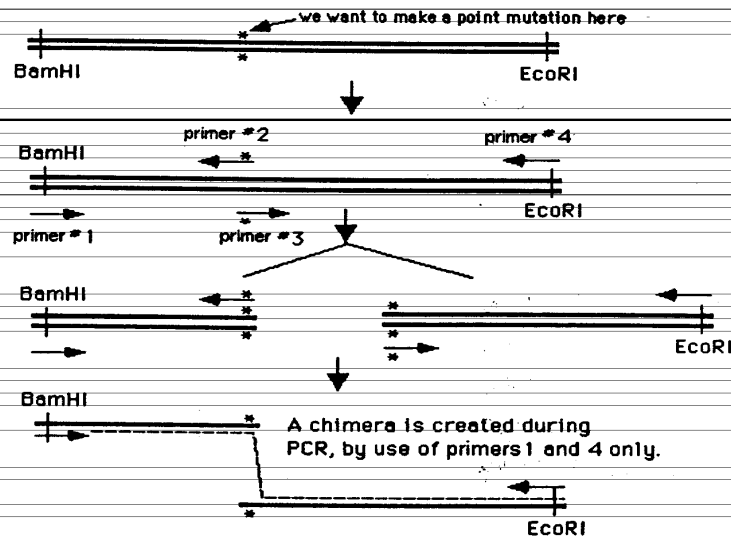
ويلاحظ أننا في النهاية قد حصلنا على قطعة جديدة من الـ DNA تم تحويل نهايتها
عن طريق إدخال تتابعات من النيوكليوتيدات غير المتوافقة mismatch (طفرات
إستبدالية) عن طريق تحويل البادئات.

٢.٢.٢. إستبدالات substitutions عند وسط الجزيء.

عادة الاستبدالات (الطفرات الموضعية) يرغب في تواجدها في وسط التتابع
(الأسكات فعل ORF ما) ونادرا عند نهاياته، ومن أجل هذا نلجأ لإستعمال عدد أكبر من
البادئات، كما هو موضح بشكل (٢-٢). فإذا ما رغبتنا في إدخال طفرة بموقع متوسط
والمثلة في الرسم بالعلامة "*" نستخدم ٤ بادئات مختلفة منها اثنين محيطة بالقطعة
وهما primer#1 و primer#4 والاثنين الآخرين يغطيان منطقة التغير (الطفور) وهما
primer#2 و primer#3، وبتعظيم قطع الـ DNA هذه سيتولد لنا قدر كبير من القطع كل
يمثل نصف القطعة الأولية parent molecule والتي تتميز أطرافها بمواقع لتعارف بعض
إنزيمات القصر (لتسهيل عمليات الكلوثة cloning فيما بعد). وفي الخطوة التالية تستعمل

القطع الجديدة في دورة جديدة من التعظيم بواسطة الـ PCR ولكن في هذه المرة سنستعمل البادئات primer#1 و primer#4 فقط وفي هذه الحالة يمكن أن يتولد لدينا قطع تمثل خليط من النصفين فيما يعرف بأسم chimera molecule حيث سيمثل الجزء الأصلي بعد إدخال بعض الإستبدالات التي تؤدي إلى عدم توافق mismatch وذلك في منتصفه (كما هو مبين بشكل: ٢-٢).

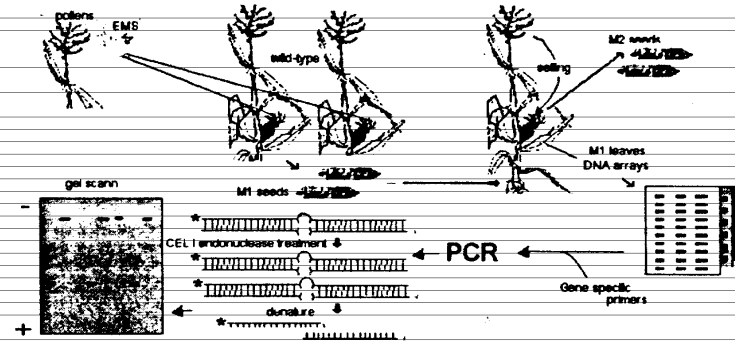
عادة "تكون" مثل هذه القطع الجديدة (الطافرة) في بلازميدات وتستغل في تجارب للتهجين بغرض الحصول على أفراد عبورية recombinants بها التتابعات الطافرة وخلفية من التراكيب البرية wild-type.



شكل (٢-٢): رسم توضيحي يمثل خطوات إستحداث طفرة بمنتصف قطعة من الـ DNA.

٢.٢.٢. استعمال المطفرات الكيميائية مع تقنيات الـ PCR.

الطرق السابقة لإستحداث طفرات موجهة الموضع site-directed عادة يمكن تطبيقها بسهولة نسبية في الكائنات ذات الجينومات الصغيرة أو المعروفة بالتحديد، أما الكائنات كبيرة الجينوم مثل نبات الذرة maize plants فكان تطبيق تلك الطرق السابقة بها غاية في الصعوبة حتى اقترح Till وزملاءه سنة ٢٠٠٤ طريقة جديدة تجمع مع استعمال المطفرات الكيميائية (مثل EMS) وتقنيات الـ PCR وتعرف هذه الطريقة بأسم TILLING وهو اختصار لـ Targeting Induced Local Lesions IN Genomes. وفي هذه الطريقة، تجمع حبوب اللقاح pollens من أحد نباتات الذرة التي تمثل صف line ما وتعامل بالمطفر القوي EMS (ethylmethansulfonate) ثم تلقح كيزان النباتات الطبيعية (والتي لها نفس الخلفية الوراثية للنباتات المعاملة) بحبوب اللقاح المعاملة، ثم زرع الحبوب نبتات الذرة التي تمثل صف line ما وتعامل بالمطفر القوي EMS (ethylmethansulfonate) ثم تلقح كيزان النباتات الطبيعية (والتي لها نفس الخلفية الوراثية للنباتات المعاملة) بحبوب اللقاح المعاملة، ثم تزرع الحبوب الناتجة للحصول على نباتات الجيل الأول الطافر M₁-generation. يمكن بالتلقيح الذاتي الحصول على نباتات الجيل الثاني الطافر M₂-generation وهكذا. نباتات الجيل الأول M₁ تعتبر خليطة وراثيا heterozygous لأي طفرة مستحدثة، ولتحديد الطفرات تجمع عينات من أوراق الجيل الأول أو صفوف الجيل الثاني حيث يستخلص الـ DNA من أنسجة الورقة لكل نبات على حدة، ثم تستعمل بادئات primers خاصة بالجين المراد إدخال طفرات عليه، ثم تعظم بواسطة الـ PCR. تجمع العينات المعظمة وتفصل الخيوط denatured عن بعضها ثم يعاد الارتباط بينها مرة أخرى re-annealed لتكوين جزيئات خليطة من الخيوط طافرة وخيوط طبيعية heteroduplexes، ثم تعامل هذه الجزيئات الخليطة بأنزيم القصر CEL I، ثم تفحص القطع الناتجة بالتفريد الكهربى حيث يمكن تحديد الصفوف الطافرة من غيرها من حجم القطع المفردة، كما هو موضح بشكل (٣-٢).

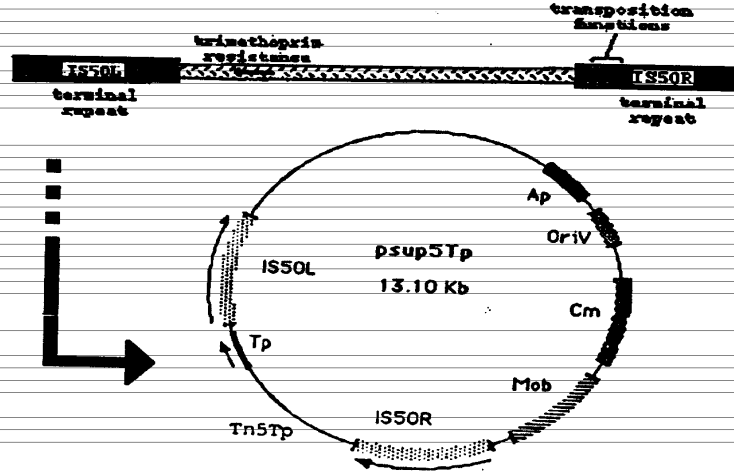


شكل (٢-٣) : رسم توضيحي لخطوات عزل الطفرات بنبات الذرة باستعمال طريقة TILLING.

٤.٢.٣. التحطّر بالإقحام Insertional mutagenesis

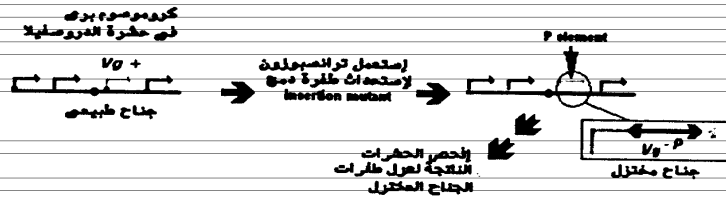
بالإضافة للطرق السابقة لإستحداث طفرات محددة الموقع تستعمل العوامل الوراثية المتحركة mobile genetic elements، خصوصا الترانسبوزونات Transposons، لقدرتها على الإندماج insertion مع DNA خلية العائل host cell الذي قد يكون كروموسوما أو فيروسا أو بلازميدا أو غيرها، وبإندماجها تؤدي إلى اضطراب الجين عند موضع الدمج مسببة فقدده لقدرته على العمل فيما يندرج تحت مجموعة طفرات الإطاحة - knockout mutants. وقد اتبعت هذه الطريقة في بادئ الأمر في البكتريا بنجاح حيث استعمل الترانسبوزون المعدل Tn 5Tp على نطاق واسع، حيث أدخل جين المقاومة للمضاد الحيوى trimethoprim محل الجين الخاص بمقاومة المضاد kanamycine في الطرز البرية، كما هو مبين بالشكل (٢-٤). وعادة ما يولف هذا الترانسبوزون مع بلازميدات معروف بها عدد من جينات مقاومة المضادات الحيوية مثل الـ Amp = Ampicilline و الـ Cm = chloramphenicol، بالإضافة لجين مقاومة المضاد Tp = trimethoprim الموجود على الترانسبوزون Tn5Tp نفسه. وعادة يتم التزاوج بين البكتريا *E. coli* الحاملة لهذا البلازميد المركب مع البكتريا المراد إستحداث الطفرات بها ثم يفحص النسل الناتج لتحديد الخلايا المقاومة للمضاد Tp كدليل على إندماج

الترانسبوزون بها ثم تفحص هذه المجموعة للبحث عن الطفرات المستحدثة أو المرغوبة. ولكن من المصاعب التي واجهاتها هذه الطريقة كانت قدرة الترانسبوزون على التنقل بفعل إنزيم الـ transposase مما يؤدي لارتداد الطفرات وفقدانها بسهولة.



شكل (٤-٢) : رسم توضيحي يمثل الترانسبوزون Tn5Tp منفردا، وبعد ربطه بالبلازميد المركب pSUP5Tp (حيث الجينات: Amp - مقاومة الأمبسلين و Cm - مقاومة الكلورامفينيكول و Tp - مقاومة تريموثوبريم و OriV - مبدأ التضاعف و mob - مبدأ النقل و الترانسبوزون Tn5Tp).

فيما بعد عرفت العديد من الترانسبوزونات في الحيوانات خصوصا الدروسفيليا والفئران، وقد استعمل العامل الوراثي المتنقل المعروف بأسم P element على نطاق واسع لإستحداث طفرات بالنسيج الجنسى لحشرة الدروسفيليا بناء على الخطوات المبينة بشكل (٥-٢).



شكل (٥-٢) : رسم تخطيطي لإستحداث طفرات لصفات الجناح في الدوسفيلا باستخدام العامل الوراثي المتنقل P.

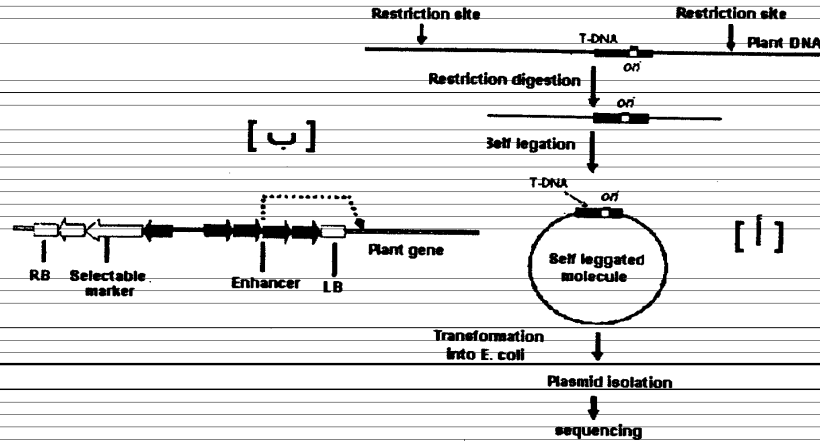
لكن كان يعيب هذه الطريقة هي الأخرى عدم الثبات بسبب تنقل

الترانسبوزونات بالإضافة إلى أن العامل P في الدوسفيلا يندمج بمواقع محددة على الكروموسوم مما يشك في عدم حيادية إستحداث مثل تلك الطفرات. ولقد أستخدم في الدوسفيلا ترانسبوزونات أخرى لتجنب تخصصية مواقع الدمج مثل الترانسبوزونات المعروفة بأسم Hermes and piggyback.

أما في النباتات الراقية خصوصا تلك التي تم الإنتهاء من التعرف على

جينوماتها الكاملة مثل نبات الـ Arabidopsis ونبات الأرز فتستحدث بها طفرات الألفام هذه بواسطة قدرة الاندماج لقطع الـ T-DNA المحورة والخاصة بالبكتيريا Agrobacterium tumefaciens المستعملة على نطاق واسع في تجارب التوليف الوراثي بالنباتات. و نبات الـ Arabidopsis هذا يعتبر نموذجا لمثل تلك التجارب، فقد استعمل على نطاق واسع لهذا الغرض لعرفتنا المستفيضة عن خرائط تنابعاته والواسمات الجزيئية العديدة المتوفرة بجينومه، وقد أستعمل الـ Arabidopsis من خلال عدة تقنيات للتطفر تؤدي إلى فقد عمل الجين loose of function عن طريق الألفام والتي تعرف بأسم التبديل بواسطة الـ T-DNA أي T-DNA tagging. أن إستعمال قدرة الاقتحام أو الاندماج التلقائية للـ T-DNA في تنابعات جينوم الـ Arabidopsis الصغير وبواسطة

الكلونة وتقنيات الـ PCR يمكننا معرفة التتابعات المحيطة بجانب منطقة الأقدام، كذلك يمكننا التنبؤ بالجين الذى فقد نشاطه نتيجة للأقدام. ولكن يجب التأكيد على أن جملة الطفرات الظاهرية المشاهدة فى مثل تلك التجارب لا يمكن إرجاعها جميعا لعملية الأقدام ذاتها لذلك وجب التوصل لطريقة لتأكيد الارتباط بين الطفرة المشاهدة وعملية الأقدام. شكل (٦-٣) يوضح أحد الاستراتيجيات المستعملة لذلك والتي تعرف باسم plasmid rescue حيث يدخل تتابع منشأ التضاعف *ori* على تتابعات الـ T-DNA حتى يمكنه الدخول والتضاعف فى البكتريا *E. coli*.



شكل (٦-٣): [١] رسم تخطيطى لخطوات طريقة الـ Plasmid rescue فى نبات الـ *Arabidopsis*. [ب] رسم توضيحي لـ T-DNA المدخل على تتابعات المستحث enhancer لتصيد الجينات.

ونتيجة للجهود الكبيرة في هذا المجال خلال السنوات القليلة الماضية، يمكننا الآن الحصول على عدد كبير من سلالات الـ *Arabidopsis* الطافرة بالأقحام insertion-mutated lines والتي يبلغ عددها حوالي ١٧٥٠٠٠ سلالة مختلفة، يمكن الحصول عليها من مراكز حفظ الأصول العالمية *Arabidopsis stock centers* كما في جامعة أوهايو بالولايات المتحدة أو جامعة نوتنجهام ببريطانيا. وهناك تقنية معاكسة ولكنها تعتمد على إستحداث طفرات بالأقحام ولكن هذه الطفرات ترجع لظاهرة إسترجاع النشاط الجيني gain of function وليس فقده. وتعتمد هذه الطريقة المعروفة بأسم "مصيدة الجينات" gene trapping اعتمادا على تراكيب مولفة من الـ T-DNA مدخل بها تتابعات من المحنات للنسخ من بروموتور قوي مثل البروموتر CaMV 35S، كما هو مبين بشكل (٢-٦). فمثل هذا الأقحام قد يستحدث طفرات مظهرية جديدة لأحد أفراد أحد العائلات الجينية gene family والتي كانت غير عاملة.

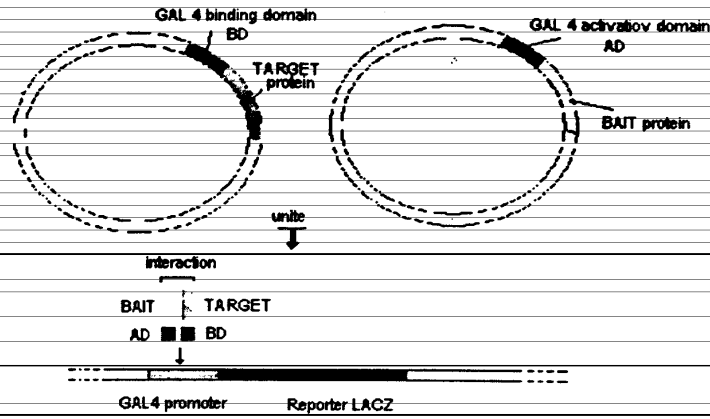
٢.٢.٢. الفعل الجيني الفائق Gene Overexpression.

هذه التقنية مثلها مثل سابقتها صممت للدراسة وظائف الجينات وهي تعتبر معاكسة للمفاهيم التي تناولناها في التقنيات السابقة للإطاحة بفعل الجين فهي على العكس تهتم بدراسة تأثير النشاط الجيني الزائد على الشكل الظهري. وفي هذه الطريقة، يتم إدخال عدد كبير من نسخ الجين تحت الدراسة إلى بلازميد أو أداة من أدوات النقل الجيني multicopy cloning vector، كما هو موضح بشكل (٢-٧). يتم بعد ذلك إدخال هذه الجينات في جينومات الكائنات تحت الدراسة، ومنها القتران للدراسة النشاط الجيني الزائد على الشكل الظاهري.

وتستعمل هذه التقنية كثيرا لدراسة العلاقة بين الجينات ومرض السرطان cancer في الثدييات.

٤.٣. دراسة التداخلات بين البروتينات .Proteins Interactions

بينها يستعمل التكنيك المعروف باسم yeast dihybrid حيث يستعمل بلازميدين، فعلى سبيل المثال الجين المنظم للنسخ المعروف باسم Gal 4 يتكون من وحدتين بروتينيتين two domains يشفر لكل منهما تتابعات محددة على الجين، أحدهما خاصة بالتنشيط activation والثانية بالارتباط بتتابعات خاصة على الـ DNA ولا بد من ارتباط الوحدتين مع بعضهما بالتوازي juxtaposition لكي يستطع القيام بعمله. ففي البلازميد الأول "يكون" الجزء الأول للجين Gal 4 والمستول عن تتابعات التنشيط ويربط معه تتابع الجين المستول عن البروتين الأول المراد دراسته bait protein، أما البلازميد الثاني فيكون به الجزء الثاني من الجين Gal 4 والمختص بالارتباط بالـ DNA مع البروتين الثاني تحت الدراسة target protein، كما هو موضح بشكل (٢-٨). ودمج البلازميدين مع بعضهما فإذا ما كان كلا البروتينين تحت الدراسة يمكنهما الارتباط مع بعضهما، فإنهما بالتالي سيسمحان بالارتباط بين وحدتي البروتين Gal 4 وبالتالي يكون قادر على تنشيط البروموتر لنسخ جين استدلال reporter gene وهو في هذه الحالة Lac Z فإذا ما أظهر فعلة دل ذلك على أن البروتينين تحت الدراسة يتداخلان مع بعضهما البعض.



شكل (٣-٨) ، رسم توضيحي لطريقة الـ Yeast dihybrid method .

٥. التعارض للإسكات الجيني بواسطة جزيئات الـ RNA الصغيرة - RNAi

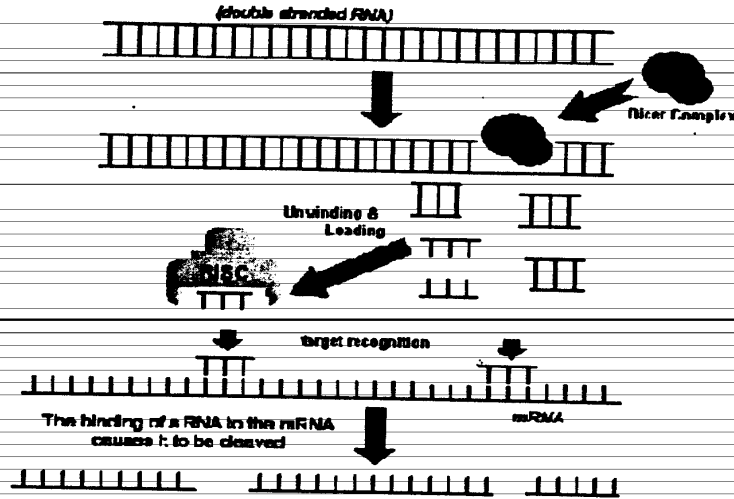
عرفت هذه الظاهر منذ وقت ليس ببعيد، عندما لوحظ في منتصف الثمانيات من القرن الماضي، أن نباتات الدخان المعدلة وراثيا transgenic plants بأحتوائها على جين الغلاف البروتيني لقيرص تبرقش أوراق الدخان tobacco mosaic virus coat protein اكتسب مناعة ضد الإصابة بالفيرس (TMV). وبعد ذلك لوحظ أن هذه الظاهرة تعتبر عامة في الكثير من الكائنات التي درست مثل الفطريات ونبات الأرابيدوسس والدودة *C. elegans* و الدروسوفيلا والثدييات، حيث عرفت بتسميات مختلفة ففي الفطريات Fungi عرفت باسم quelling وفي النباتات باسم RNAi. وخلال السنوات العشر الماضية تعرفنا على نظام مشابه يتحكم أيضا في الإسكات الجيني ولكن بواسطة جزيئات صغيرة تعرف باسم micro RNA [miRNA]. وتتفق كل تلك الظواهر في القدرة على إسكات الجينات ويتم ذلك من خلال الخطوات التالية:

- بناء جزيئات من الـ RNA مزدوج الأذرع يسمى ds RNA بواسطة الأنزيم RNA dependent RNA polymerase ويعرف اختصارا (RdRP)،

وهو قد يشابه في عمله أنزيمات بناء الفيروسات النباتية (viral replicases). ويمثله الناتج الجيني لجينات مثل *SDE-1* في الأرابيدوبسيس أو *EGO-1* في الديدان *C. elegans*.

- يقوم أحد أنزيمات الهدم الـ RNA الداخلي (endonuclease) من عائلة RNase III ، *E. coli* ، (، أطلق عليه اسم Dicer) بهدم الـ dsRNA إلى قطع صغيرة حجمها يتراوح بين ٢١-٢٥ نيوكلويدة، لها امتداد من الطرف 3'-OH نيوكلويدتين (3' n: overhang in each sense and antisense). وهذا الأنزيم يحتاج في عمله إلى طاقة ATP وقطع مزدوجة من الـ RNA لا يقل طولها عن 39 bp ، وقد يتماثل هذا الأنزيم مع الناتج الجيني *RDE-4* الخاص بـ *C. elegans*.
- يقوم أنزيم من نوع الـ Helicase (غير معروف خواصه على وجه الدقة) بفك ارتباط ذراعي الـ ds RNA وتحويله إلى قطع مفردة من الـ short interfering RNA (si RNA). ويعرف هذا المعقد الإنزيمي اختصاراً بـ RISC للمصطلح RNA-induced silencing complex.
- تتعرف قطع الـ siRNA مع التتابعات المناظرة لها على الرسائل المستهدف والمراد إسكاته (target mRNA)، وعن طريق البلمرة polymerization أو اللصق ligation تتكون جزيئات مزدوجة يتم هدمها بسرعة بواسطة الـ Dicer.

أما في حالة الدروسوفيلا والشد ييات ومنها الإنسان، فإن ميكانيكية حدوث الأسكات الجيني بتداخل الـ RNA مازالت غير واضحة. ولكن بينت التجارب أن ارتباط الـ si RNA مع الرسائل المستهدف لا يحتاج لأنزيم RdRP وان النهايات 3'-OH لهذه القطع الصغيرة من الـ si RNA لا تعمل كبادئ للبناء (primers) كما في الأنظمة الأخرى. والشكل (٢-٩) يلخص هذه الخطوات.

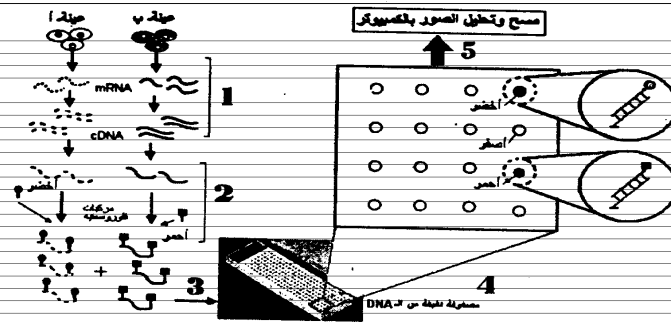


شكل (٢-٩) : رسم توضيحي لخطوات الإسكات الجيني بواسطة RNAi.

٦.٢. المصفوفات الجينومية Genomic Micro-arrays.

أصبح من الواضح لدينا أن الكائنات على اختلاف أنواعها وأشكالها تحتوى على أعداد كبيرة من الجينات، حيث يتم التعبير عن نسبة عالية منها فى الزمان والمكان المحدد من حياة الكائن. هذا القدر من العمل يستوجب درجة عالية من التنسيق فيما بينها، ولكن وللأسف فإن الأساليب التقنية التى كانت متاحة من قبل كانت لا تسمح إلا بدراسة التعبير الجينى لكل جين على حدة وبالتوالى، حيث كانت الصورة الشاملة للتعبير الجينى غير متاحة. ولكن فى السنوات العشر الماضية ظهرت تقنية جديدة تتيح دراسة التعبير الجينى لمجموعة من الجينات (قد تصل لعدة آلاف فى المرة الواحدة). هذه التقنية الجديدة عرفت بعدة مسميات، فأطلق عليها اسم مصفوفات الـ DNA الدقيقة DNA microarray وسميت بعد ذلك "رقائق الـ DNA" DNA chips لكونها دقيقة جداً وتشابه تكنولوجيا صناعة الدوائر المدمجة integrate circuits فى الإلكترونيات مع التأكيد على أنها لا تمت لعلوم الإلكترونيات بأى صلة. وتفصيلاً يمكن القول، بوجود

نوعين من تلك المصفوفات، الأول هي المصفوفات الكبيرة (نسبيا) DNA macroarrays وهي التي تتكون من نقاط من عينات الـ DNA تزيد عن ٥٠٠ ميكرون لكل نقطة - في هذه الحالة يمكن قراءة الصور بماسحات scanners عادية، والمصفوفات الدقيقة DNA microarrays حيث لا يزيد حجم العينة عن ٢٠٠ ميكرون وهنا يجب قراءة الصور الناتجة بماسحات خاصة. ومن التسميات الشائعة أيضا - الرقائق الجينية gene chips وقد استعمل هذا المصطلح في الكثير من البحوث المنشورة في الدوريات العلمية المتخصصة، إلا أن هذا الاسم يعتبر ملكية أدبية وتجارية لشركة Affymetrix Inc. لتسجيلها منتجا تجاريا بهذا الاسم، مما دفعها للتقدم بعدة دعاوى قضائية ضد من يستعمل هذه التسمية دون إذن مسبق منها. وبناء على هذا فقد اخترت اسم الرقائق الجينومية genomic chips لتسمية هذه التقنية دون غيرها من المسميات. شكل (١٠-٢) يوضح أهم الخطوات لتنفيذ هذه التقنية.



شكل (١٠-٢) : رسم توضيحي للخطوات الرئيسية لدراسة التعبير الجيني من خلال استعمال المصفوفات الجينومية.

يمكننا أن نلخص الخطوات المتبعة في تقنية المصفوفات الجينومية في خمس أقسام رئيسية والمرقمة في شكل (١٠-٢)، وهي كالتالي :

١. تحضير عينات الـ DNA تحت الدراسة، أو ما يسمى بالـ DNA المستهدف -

"target- DNA" = عادة يستخلص الـ RNA الكلى من عينات مختلفة تمثل

ظروفا تجريبية مختلفة ومنها تعزل المرسلات mRNAs كيميائيا، وعن طريق

النسخ العكسي reverse transcription يتم بناء جزيئات الـ cDNA الكاملة،

والتي تمثل المعين النسخي transcriptomea لهذه العينات.

٢. تعليم labeling للجزيئات المستهدفة بواسطة أصباغ فلوروسنتية

fluorescent dyes وعادة تستعمل مشتقات السيانين مثل Cy3 (خضراء اللون)

و Cy5 (حمراء اللون).

٣. تنقل الجزيئات المستهدفة والمعلمة بالأصباغ الفلوروسنتية لتهجينها

hybridization مع أحد المصفوفات الجينومية بما تحمله من جزيئات الـ DNA

العيارية أو القياسية أو ما يعرف باسم probe DNAs، عن طريق تتابع الدورات

المتحكم فيها من الـ denaturation/annealing بأستعمال المحاليل ودرجات

الحرارة المناسبة. والفحص الدقيق يوضح أن الأقرانات بين الجزيئات المستهدفة

(الملونة إما أحمر أو أخضر) سيتم مع جزيئات الـ probe الكاملة لها وغير المعلمة،

وعليه في حالة ارتباط مع المصفوفة فإن اللون الأخضر سيحدد جينات العينة

الأولى بينما اللون الأحمر سيحدد جينات العينة الثانية أما اللون الأصفر

فسيجدد أماكن الارتباط المشترك لجينات العينة الأولى والثانية معا، أما إنعدام

اللون فسيمثل عدم وجود تهجين (نشاط جيني)، كما هو مبين بالرسم الكبير

(القصى اليمين) في شكل (٣-١٠).

٤. تبدو هذه الخطوة وليست في سياق التتابع لهذه التقنية وهذا صحيح فيجب أن

تكون سابقة لبداية العمل، ومع هذا يجب التعريف بها - ألا وهي صناعة

المصفوفة ذاتها. تصنع المصفوفات أو الرقائق من الزجاج عادة أو من مواد

بلاستيكية شفافة في حالات نادرة بعد تغطيتها بطبقة رقيقة من مادة

poly-L-Lysine وتوضع عليها بترتيب دقيق عينات الـ DNA العياري probes،

بأستعمال طرق أوتوماتيكية robotic لوضع (طبع) قدر ضئيل ومحدد من هذه

الـ probes على هيئة نقاط وهذه الطريقة مماثلة لطبع الدوائر الإلكترونية بما

يسمى photolithography. ويتراوح عدد هذه الـ probes بين المئة حتى عشرات الآلاف، قد يمثل كل منها تتابعات صغيرة أو cDNA (جين) أو كروموسوم ونأمل أن تمثل الرقائق في المستقبل القريب جينوما بأكمله.

٥. الخطوة الأخيرة في هذه التقنية غاية في الأهمية - وهي مسح وتحليل الصور

الناجمة من تلك التجارب وهي تحتاج إلى بعض التفصيل.

١.٦.٢. مسح وتحليل المصفوفات Scanning & Analysis of Arrays

عادة يتبادر إلى أذهاننا في الدراسات المعتمدة على استعمال المصفوفات

الجينومية - سؤالين - الأول يتسائل عن شدة التعبير لكل جين من جينات هذه المصفوفة

تحت الدراسة - والتساؤل الثاني يتمحور حول تقدير العلاقات بين قدر التعبير الجيني

لجينات المصفوفة تحت الدراسة. الأجابة على السؤال الأول عادة سهلة وممكنة، فبتوفر

الأجهزة الحديثة التي تعتمد على تكنولوجيا الليزر يمكن إستحساس excitation وتقدير

قدر الامتصاص الضوئي optical absorption بدقة متناهية عند أطوال الموجات للأصباغ

الفلوروسنتية المستخدمة وربطها بقدر النشاط الجيني لكل منها. أما الأجابة على السؤال

الثاني فهي التي ستحتاج منا جهدا وتديقا أكبر من خلال إستعمال الأساليب الحسابية

والإحصائية، وغالبا ما تكون غير فاطعة بل ستكون مقبول فقط عند درجة من درجات

الاحتمال.

في مبدأ الأمر اتبعت طرق حساب " التجميع العنقودي " للبيانات

cluster analysis لتقدير العلاقات بين مجاميع الجينات النشطة تحت ظروف التجربة.

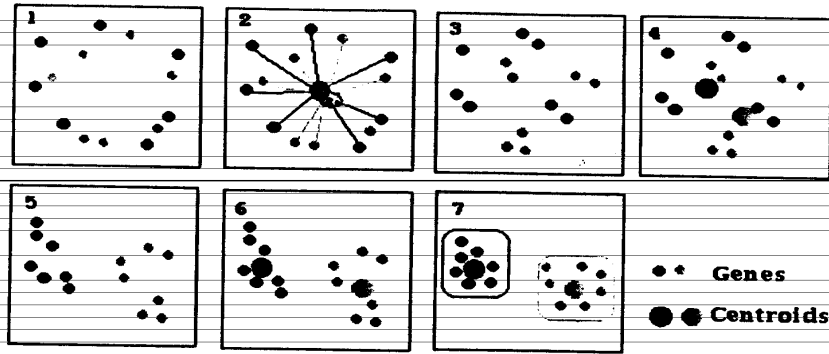
فمراجع الكمبيوتر وطرق حساب هذا التجمع لها عدة طرز ولكن أكثرها شيوعا وأبسطها

هي طريقة k-means clustering حيث يتم أولا إعطاء كل جين في المجاميع النشطة

قيمة رقمية ثم تحدد المتوسطات (مركز الثقل - centroid) لكل مجموعة ثم خطوة

خطوة تجمع كل مجموعة متشابهة التعبير مع بعضها البعض وصولا للفصل ما بينهم ما

يمكن. شكل (٢-١١) يمثل رسوما توضيحية لخطوات طريقة k-means clustering.



شكل (١١-٢) : رسوم توضيحية لبعض من خطوات طريقة التجميع العنقودي *k-means*

clustering لتقدير درجات التشابه في التعبير الجيني لجمعتين من الجينات النشطة المميزة

بالنقاط الصغيرة الرمادية والسوداء، بينما مركز النخل لكل مجموعة ممثل بالنقاط الكبيرة.

طرق التجميع العنقودي على اختلاف أشكالها وصورها لاشك لها أهمية تاريخية

لكن الآن لا يفضلها عدد كبير من الدارسين بل يتجنبوها لعدم دقتها ويفضلون اللجوء إلى

لطرق الإحصائية التقليدية للمقارنة بين مجاميع البيانات مثل اختبار t ، كما هو مبين

بالمعادلة التالية:

$$t = (x_{i,1} - x_{i,2}) / ((s_i + s_0)/2).$$

حيث $x_{i,1}$ هو متوسط تأثير الجين i في المجموعة ١ و $x_{i,2}$ متوسطه في المجموعة s_i و

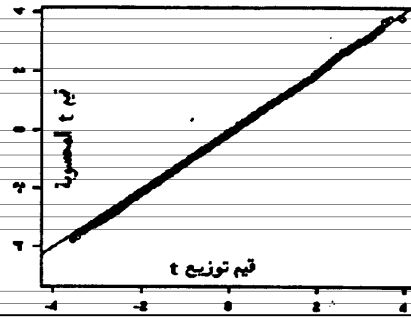
لخطأ القياسى داخل المجاميع و s_0 متوسط توزيع الخطأ القياسى في كل المجاميع. عادة

يكون اختبار t في حد ذاته كافيا لذلك وجب اختبار قيم t المحسوبة منسوبة لتوزيع

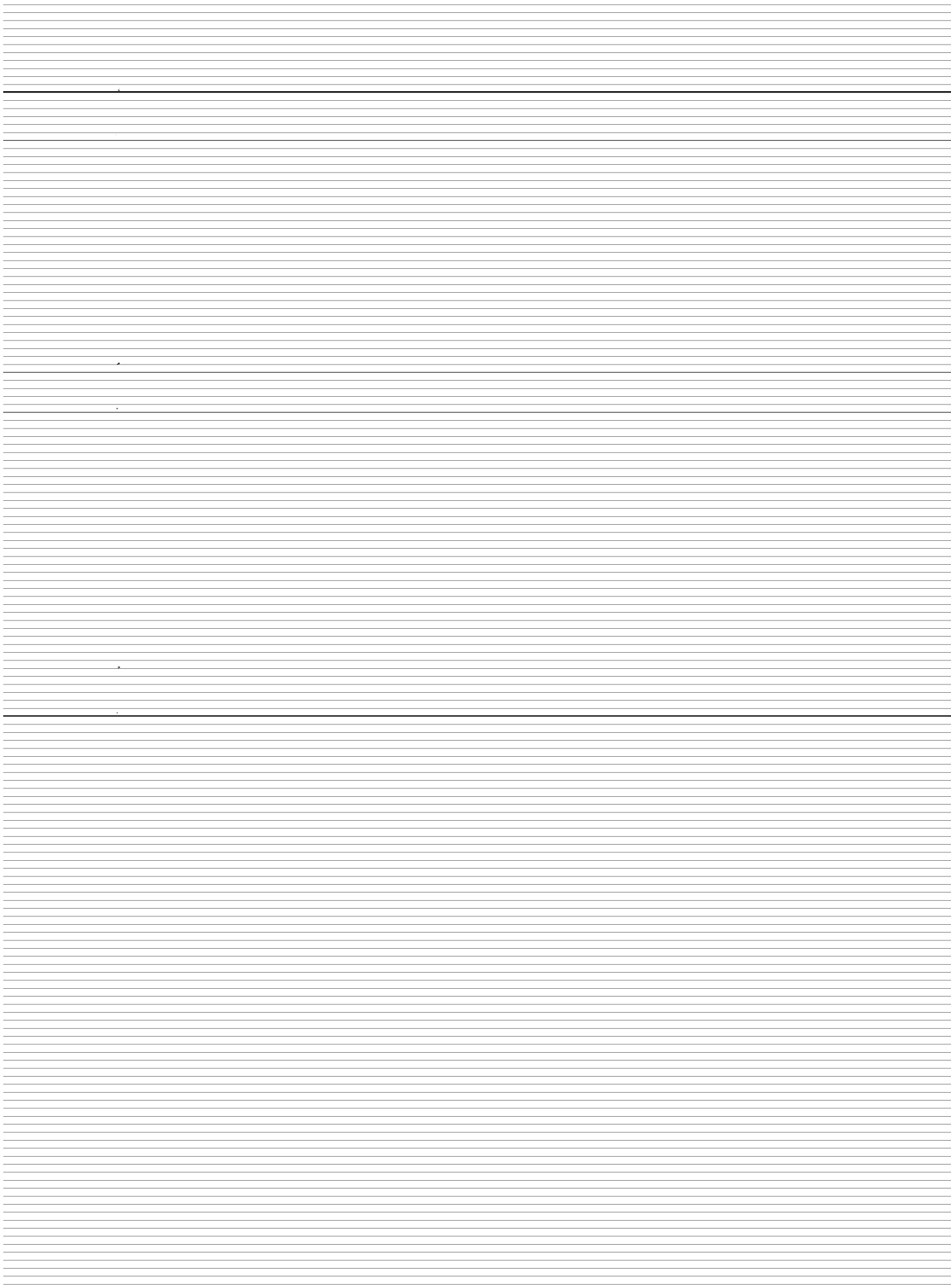
قيم t ، حيث عادة يلاحظ أن قيم t لا تظهر اختلافات بين الجينات في المجاميع وأنها

تتبع توزيع منتظم ممثل بخط مائل بدرجة ٤٥، أما الجينات المختلفة فستنحرف عن هذا

التوزيع بدرجات متفاوتة، كما هو مبين بالعلاقة البيانية التالية:



ومع ان اختبار توزيع t يعتبر مقبولا بشكل اولي الا ان الطرق الاحصائية المستعملة
 لدراسة الاختلافات في تجارب المصفوفات الجينومية اصبحت اليوم أكثر تعقيدا ودقة،
 وبعض من تلك الطرق الاحصائية الشائعة سنتناولها بقدر من التفصيل في الباب القادم.



٤. الأساس الحسابي للمعلوماتية الحيوية

Mathematical Bases of Bioinformatics

إعداد: أحمد المتينى

يمكن القول بأن المعلوماتية الحيوية هي التطبيق المباشر لتكنولوجيا المعلومات (information technology) في مجال البيولوجيا الجزيئية (molecular biology) حيث تتم الاستعانة بأساسيات علوم الرياضيات التطبيقية (applied math) والكمبيوتر (computer) والإحصاء (statistics) لتخزين وتبويب وتفهم المعلومات البيولوجية. فكم المعلومات المنتجة يوميا في مجال البيولوجيا الجزيئية والجينومكس فاق الحدود المتوقعة وأصبح الإلمام بها أو حتى ببعضها منها بالجهود الذاتية من شبه المستحيل، فعلى سبيل المثال قدرت الأوراق العلمية المنشورة في هذا المجال في خلال السنوات القليلة الماضية بنحو ٢٠٪ من جملة البحوث المنشورة في مجال الطب والبيولوجيا عموما وذلك حسب إحصائيات قاعدة بيانات المكتبة الأمريكية القومية للطب (NLM) التابعة للمعهد القومي للصحة (NIH) والمعروفة اختصارا باسم "PubMed" [www.nlm.nih.gov]. كذلك نشرت صحيفة الإيكونست في سنة ١٩٩٩ مقالة لـ Anthony Kervalege، أحد كبار يبحاث شركة Celera الشهيرة، قدر فيه حجم المعلومات التي يمكن أن ينتجها معمل مختص بهذا المجال بحوالي 100 Gb يوميا!

إن كمية المعلومات المنتجة في هذا المجال ضخمة للغاية ومهولة لذلك وجب الاستعانة بالكمبيوتر لإنشاء قواعد للمعلومات (databases) تبويب وتخزين بها تلك المعلومات ثم تستغل برامج الكمبيوتر (software) للبحث وتحليل وإستخلاص المعلومات، سوف نتناول بعض من هذه المواضيع بشكل أوسع في الأبواب القادمة. أن نتاولنا لهذه المعلومات من خلال مسح قواعد المعلومات المتاحة واسع ومتعدد حيث يتميز بالديناميكية العالية مقارنة بغيره من فروع العلم، لكن يمكننا أن نلخص أهم الأهداف في النقاط التالية:

- **البحث الإستخلاصي (التجميعي) Annotation search** ، حيث تستعمل مفاتيح لتصفح قواعد المعلومات باستعمال كلمات فاتحة (keywords) أو أسماء العلماء (authors) أو المواضيع المرتبطة، بفرض زيادة أو تجميع معلومات عن موضوع ما مثل معرفة — ماهو الجين (الجينات) المسئول عن مرض cystic fibrosis في الإنسان؟
- **البحث عن المثل أو النظم Homology (similarity) search** ، حيث يتم حصر وتصفح قواعد المعلومات المتاحة بفرض التعرف على نماذج متماثلة لما لدينا من معلومات مثل معرفة — هل تتابعات الجين المعرف لدينا تناظرها أو تماثلها تتابعات في كائن آخر؟
- **البحث عن المميزات Pattern search** ، حيث يتم حصر قواعد المعلومات بحثاً عن أجزاء مميزة للجزيئات قد تستعمل في معرفة وظائفها مثل — ما هي التتابعات المميزة لجين ما تعمل كمناطق ارتباط عوامل البنى (Irritation) أو الحس (enhancing) أو الإسكات (silencing) خلال النسخ (transcription) الجيني.
- **التنبؤ Prediction** ، استغلال المعلومات المتاحة بقواعد المعلومات للتوصل لإستنتاجات جديدة مثل التنبؤ بالتراكيب الثانوية (secondary structure) لبروتين ما من بيانات تركيبه الأولى (primary structure).
- **المقارنات Comparisons** ، استغلال قواعد المعلومات المتاحة لمقارنة تتابعات معلومة من (الأحماض النووية أو البروتينات) بغيرها من التتابعات مثل البحث عن العائلات الجينية (gene families) في جينوم ما.

خلال عمليات الفحص والتنقيب والإستخلاص لقواعد المعلومات هذه لابد لنا من أن نعتد على مبادئ علم الإحصاء وبرامج و حسابات (algorithem) الكمبيوتر، وفيما يلي سنحاول الألام ببعض من تلك المواضيع ذات الارتباط المباشر بمجال المعلوماتية الحيوية دون توسع و الذي ليس مجاله هنا.

١.٤. نظرية الاحتمالات Probability theory.

تلعب الاحتمالات دورا هاما في تفسير سلوك الكائنات الحية (كما أكدنا سابقا)، ويقيد الاحتمال هو ببساطة الإمكان المقرون ببعض الثقة وأيضا ببعض الشك في وقوع حدث (event) معين. فقد تكون الثقة معدومة في وقوع حدث ما ومع ذلك قد يفترض أنه محتمل الوقوع، ومن ناحية أخرى قد تكون الثقة في حدوث حدث تصل لحد التأكد ومع ذلك يكتفى بذكر أن الحدث محتمل الوقوع. وعلى ذلك يشمل هذا المعنى العام للإحتمالات كل من الاحتمالات الضعيفة والقوية على حد سواء. وقد يكون الحدث بسيط (simple) أو مركبا (compound)، أي يتكون من جملة أحداث بسيطة.

الاحتمال الرياضي للحدث البسيط هو عبارة عن عدد مرات نجاحه أو فشله منسوبة للعدد الكلي لعدد مرات النجاح وال فشل. وهو مقياس موضوعي مستقل تماما عن الهوية الشخصى. فإذا ما رمزنا للنجاح بـ a ولل فشل بـ b فإن احتمال نجاحه $a/a+b = (p)$ و أن احتمال فشله $b/a+b = (q)$ ، ومجموع احتمالات النجاح والفشل تساوى الواحد الصحيح $p + q = 1$. ومن هذه المعادلة يمكن إستخلاص الأساسيات التالية: $p = (1 - q)$ و $q = (1 - p)$

وفي حالة الأحداث المركبة، إذا كانت الأحداث البسيطة المكونة لها غير مستقلة عن بعضها، بمعنى إذا وقع أحدها منع وقوع الحدث أو الأحداث الأخرى في هذه الحالة تجمع احتمالات الأحداث الفردية، أي

$$P = p_1 + p_2 + p_3 + \dots + p_n.$$

وإذا كان الحدث المركب مكونا من أحداث بسيطة مستقلة عن بعضها البعض، أي إذا وقع حدث منها لا يؤثر في وقوع الأحداث الأخرى، في هذه الحالة تضرب احتمالات الأحداث البسيطة، أي

$$P = p_1 \times p_2 \times p_3 \times \dots \times p_n.$$

لتقدير الاحتمال الحسابى يجب أن نلم أيضا ببعض المعلومات عن التباديل (permutations) والتوافيق (combinations). التباديل = هي عدد الترتيبات المختلفة الممكن تكوينها من أحداث مختلفة عن بعضها البعض، فمثلا الحروف الثلاث a و b و c

يمكن ترتيبها في ٦ طرق مختلفة (تباديلها): abc, acb, bac, bca, cab, cba أي أن التباديل = $3 \times 2 \times 1$ ، فإذا كان لدينا عدد من الأحداث قدره n فإن التباديل = $3 \times 2 \times 1 \times \dots \times 2 \times 1 = n!$ أي يساوي مضروب n (n -factorial). وإذا ما كان لدينا ثلاث أحداث مثل a و b و c وأردنا معرفة عدد التوافيق الممكنة (إذا ما اختيرت اثنين اثنين) فستكون التوافيق الثلاث الممكنة ab, bc, ac ويمكن حساب التوافيق من المعادلة العامة:

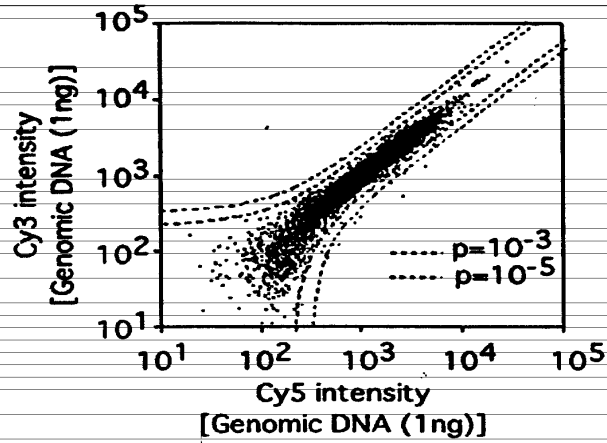
$${}^nC_r = \frac{n!}{(n-r)! \cdot r!}$$

حيث أن n = عدد الأحداث و r = عدد التوافيق المطلوبة.

١.١.٤. قواعد عامة عن توزيعات الاحتمالات Probability distributions.

فيما يلي بعض من الحقائق والقواعد الإحصائية الحسابية التي قد تكون مفيدة في الدراسات البيولوجية عموماً وفي الجينومكس والمعلوماتية الحيوية خصوصاً:

- أكثر التوزيعات استعمالاً لقياس حدود الثقة هي المنحنى الطبيعي Normal distribution والذي يعرف أيضاً باسم Gaussian's distribution وذلك للبيانات المستمرة وإستعمال مثل تلك التوزيعات يمكن الباحث من تحديد حدود الثقة لبياناته وتحديد الأخطاء التجريبية عند مستويات الإحتمال المطلوبة، كما هو مبين بشكل (١-٤) لبيانات الـ DNA microarray المتحصل عليها في إحدى الدراسات



شكل (١٤) ، توزيع بيانات الـ DNA microarray بين حدود الثقة عند احتمال $P < 0.001$ و $P < 0.00001$ باستعمال قياسات المنحنى الطبيعي.

- يستخدم توزيع الـ hypergeometric distribution في حساب الإحتمالات عند الاختيار من أحداث مستقلة لا تعوض (without replacement) ، لهذا الفرض قد تستعمل المعادلة التالية:

$$P = 1 - \sum_{x=0}^{k-1} \frac{\begin{bmatrix} k \\ x \end{bmatrix} \begin{bmatrix} N-k \\ n-x \end{bmatrix}}{\begin{bmatrix} N \\ n \end{bmatrix}} \quad \text{where} \quad \begin{bmatrix} n \\ k \end{bmatrix} = \frac{n!}{k! (n-k)!}$$

N = total population

k = number of successes in total population

n = size of sample population

x = number of successes in sample population

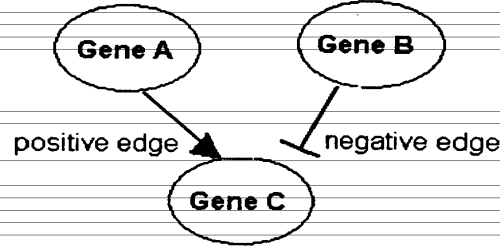
وتطبيق هذه المعادلة يمكن تمثيله بتقدير العنوية لبيانات الـ microarray لأثنين من السلالات المختلفة لنفس الكائن، كما هو موضح بشكل (٢-٤).

٢.٤. طريقة بيزا الإحصائية Bayesian statistics.

$$\ell = \frac{L(x \setminus H_0)}{L(x \setminus H_1)} \quad \& \quad \chi^2 = -2 \ln \ell = -2 \ln \left[\frac{L(x \setminus H_0)}{L(x \setminus H_1)} \right]$$

حيث H_0 تمثل الفرضية الأولى و H_1 تمثل الفرضية الثانية، وتقاس العنوية باختبار مربع كاي χ^2 .

يعتبر الوراثيون أن هذه الطريقة مناسبة لكونها بسيطة في تمثيل الجينات كذلك يمكن أن تعكس العلاقات لتأثير الجينات على بعضها البعض، والاستفادة من المعلومات السابقة في وضع الفرضيات وكذلك لقدرتها على التعامل بكفاءة مع البيانات المتداخلة مثل بيانات DNA microarrays. وقد أثبتت عن هذه الطريقة طريقة تسمى Bayesian network حيث تمثل العلاقات بصورة وصفية qualitative أولا ثم تدرس كمياً quantitative ثانياً، وهذا الأسلوب يستخدم في الوراثة بكثرة. فعلى سبيل المثال يمثل الجين بانتفاخ node ومن حوافه تخرج أسهم التأثير فإذا كان موجبا مثل براس السهم وإذا كان التأثير سلبيا مثل بالخط المنتهى كما هو موضح بشكل (٢-٤).



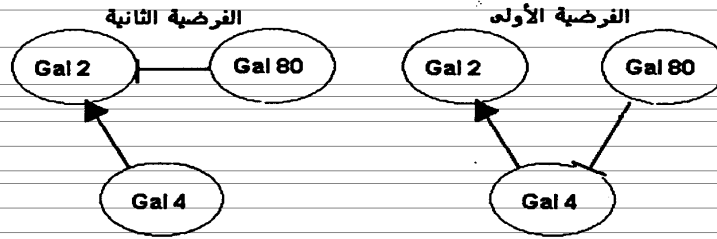
شكل (٢-٤) ، تمثيل للعلاقة بين الجينات بناء على Bayesian network.

والرسم السابق يمثل الجزء الوصفي للطريقة حيث مثلت النظرية الفرضية بناء على النتائج التجريبية، أما الشق الكمي للطريقة فيعتمد على تقدير الاحتمالات لدى تأثير النشاط الجيني لكل من الجينين A و B على نشاط الجين C سلباً أو إيجاباً، كما هو موضح بالجدول التالي:

الجين A	الجين B	احتمال التأثير الإيجابي على الجين C	احتمال التأثير السلبي على الجين C
+	+	0.6	0.4
-	+	0.01	0.99
+	-	0.99	0.01
-	-	0.4	0.6

والبيانات الموجودة بالخانة المظلمة من الجدول السابق تعنى أن زيادة نشاط up-regulation لكل من الجينين A و B سيؤدي إلى زيادة نشاط الجين C (up-regulated) بأحتمال قدره ٦٠٪ وباحتمال قدره ٤٠٪ فإن زيادة نشاط A و B سيؤدي إلى خفض نشاط الجين C (down-regulation).

كذلك تستخدم هذه الطريقة في تفهم طبيعة تأثير الجينات على بعضها البعض، ففي دراسة قام بها Hartemink *et al.* سنة ٢٠٠١ باستغلال بيانات ٥٢ محاولة للـ DNA microarrays لمجموعة من الجينات المتداخلة في التأثير على تمثيل الجلاكتوز Galactose regulatory network، وفي جزء محدد من تلك الدراسة (على سبيل المثال) اقترحوا فرضيتين لتأثير ثلاث من هذه الجينات، كما هو مبين بالشكل (٤-٤).



شكل (٤-٤) : فرضيات تأثيرات الجين Gal 80.

ولإختبار أرجحية أى من النظريتين الفرضيتين السابقتين تستخدم Bayesian network statistics عن طريق حساب معامل أو قياس يسمى Bayesian score (M)، يحسب من العلاقة الآتية:

$$M = \log [P(M/D)] = \log [P(M)] + \log [P(D/M)] + c.$$

حيث M - النموذج (الفرضية) و D - نتائج الـ microarray المتحصل عليها و C - ثابت. وبحسابات الأرجحية maximum likelihood method تحسب المعاملات لكلا الفرضيتين، فوجد أن هذا المعامل للفرضية الأولى يساوى 44.0 - بينما كان المعامل للفرضية الثانية يساوى 34.5 -، وعليه يمكن الاستنتاج بأن الفرضية الثانية أكثر أرجحية لتفسير النتائج المتحصل عليها.

ولكن أتضح أن Bayesian statistics ليست مناسبة للعديد من الحالات لذلك تستعمل طرق إحصائية أخرى مثل برامج ماركوف الإحصائية لهذا الغرض.

٢.٤. برامج (نماذج) ماركوف Markov Models .

هذا النوع من البرامج والطرق الإحصائية هو الأكثر استعمالاً في مجال المعلوماتية الحيوية لكونها أكثر ملائمة لأهداف تلك الدراسات، ولكن الأساس الرياضي الإحصائي لهذه البرامج غاية في التعقيد وتحتاج المتخصصين في الرياضة التطبيقية وعلوم الكمبيوتر للتعامل معها. هذه الحقيقة تستدعي لفت الأنظار لأهمية الفريق البحثي في هذا المجال الحديث من المعرفة، ففي مجال الجينومكس والمعلوماتية الحيوية يجب أن يتعاون الوراثيون مع المتخصصين في مجالات عديدة مثل الكيمياء الحيوية Biochemistry والفزياء الحيوية Biophysics وعلماء الرياضة Mathematicians وغيرهم لإنجاز الأهداف المرجوة. بالنسبة لنا نحن البيولوجيون لايهمنا التبحر في الأسس الرياضية لهذه البرامج، بل يهمنا تفهم الإمكانيات والأهداف العامة لتلك البرامج لتطبيقها فيما يخصنا من دراسات، وحسبنا الله أن هناك العديد من برامج الكمبيوتر software التي يمكن أن تساعدنا في هذا التطبيقات.

وفي مجال الجينومكس والمعلوماتية الحيوية تستغل برامج ماركوف الإحصائية لتحديد أرجحية التوافق matching أو عدم التوافق mismatching أو إمكانية الحذف deletion أو الأضافة insertion لتتابعات من الأحماض النووية أو البروتينات لتحديد مصادرها أو معرفة أسلافها ancestors المحتملة. وهذه البرامج الإحصائية تقلد العقل الإنساني في تعامله مع حالات عديدة مثل إلتباس الألفاظ عند أحد المتكلمين، فقد يكون ذلك لإختلاف اللهجات أو للإصابة بالتهاب في الحنجرة أو لعب خلقي في المتحدث يجعل من مخارج الألفاظ غير واضحة أو لعدم التمكن من اللغة وغيرها، ولكن عادة، المستمع يتوصل لفهم هذا الحديث حيث يقيس في عقله هذه الألفاظ المبهمة على نماذج لكلمات قريبة منها في النطق ومنها يختار الأقرب للمعنى، فمثلاً في حديث ما، قال المتحدث في سياق كلامه " المؤمن كيس فطن" وطبعاً لا يستقيم المعنى فالمقصود (كما يعيه أغلبنا دون عناء) هو " المؤمن كيس فطن".

ودون الدخول في التفاصيل الحسابية يمكن تشبيه برامج ماركوف هذه بلعبة game من ألعاب التسلية (تعتمد على الحظ الحسوب) تستعمل فيها رقعة من المربعات مثل رقعة الشطرنج أو لعبة السلم والثعبان مثلا، ولكنها لعبة غريبة ومملة ! وقواعد اللعبة يمكن أحماها في :

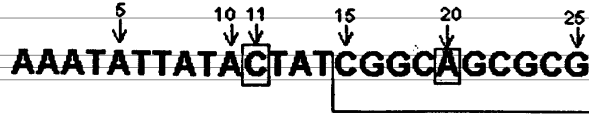
(١) كل مربع يوفر لك مجموعة من الرموز أو الحروف (في حالة الـ DNA فلدك ٤ حروف وفي البروتينات لديك ٢٠ حرفا) وتكرارات أو نسب هذه الحروف لبعضها البعض متغير (مثلا في الـ DNA الحرف A له نسبته والحرف G له نسبته كذلك الحال مع الحرفين C و T)، المهم أن كل مربع يعطى نسب مختلفة عن المربعات الأخرى (فمربع به نسب عالية من A ونسب منخفضة من الثلاث الباقية والمربع المجاور يعطى نسب عالية من G ونسب منخفضة من الثلاث الباقية وهكذا.....).

(٢) داخل كل مربع، كل حرف له قيمة حسابية تراوح بين الواحد الصحيح والصفر تميزه بناء على تكراره.

(٣) في نهاية كل دور من أدوار اللعبة يحسب مجموع الدرجات الكلية التي حصلت عليها من مجمل إنتقالاتك بين المربعات فيما قد يسمى بالنتيجة النهائية final score، والذي يكسب الدور هو من يجمع أكثر النقاط (ليس دائما).

(٤) إذا إنتقلت من مربع إلى آخر لابد أن تنال عقوبة penalty (خصما من الدرجات)، لكن هناك مربعان بالرقعة – مربع الحذف deletion ومربع الأضافة insertion – يمكنك الانتقال إليهما في أي وقت دون عقوبات لكن لن تحصل على أي درجات منهما. وإذا ما كانت هذه قواعد اللعبة، فإن الهدف من هذه اللعبة هو التداول لتتابع ما من الـ DNA لمعرفة أرجحية إنتمائه لتتابع آخر عن طريق جمع أكبر قدر من النقاط مع تجنب (بقدر المستطاع) العقوبات.

يبدو لي أن مفهوم اللعبة قد أتضح بعض الشيء، ولكن ليس بالقدر الكافي لذلك دعنا نلعب دورا مبسط في هذه اللعبة الغريبة. فإذا ما كان لدينا تتابع من ٢٥ نيوكليوتيدة من الـ DNA على النحو التالي:



وإذا كان المتاح في هذا الدور من اللعبة مربعين اثنين فقط، كما هو مبين على النحو التالي:

مربع (٢)	مربع (١)
توفر بقدر كبير كل من G's و C's ولكن يسمح بالحرفين A's و T's بندرة	توفر بقدر كبير كل من A's و T's ولكن يسمح بالحرفين G's و C's بندرة

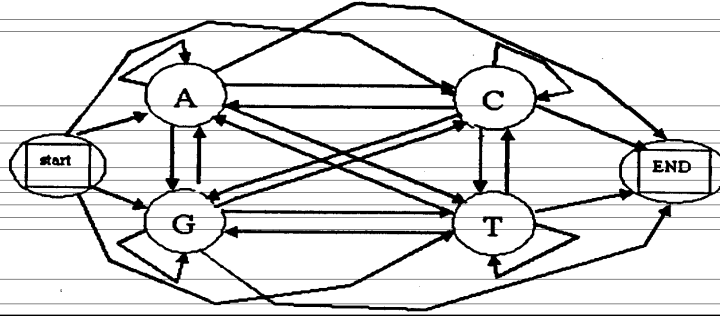
لكسب هذا الدور (بمعنى معرفة الترتيب المعطى)، يمكنك أن تبقى في المربع الأول حيث تتلقى دائما درجات عالية عن الحروف من الأول حتى الحرف رقم ١٥ فكلها من الحروف A أو T كمواصفات المربع، الاستثناء هو الحرف رقم ١١ فهو C حيث سيكون لك الخيار في أن تنتقل للمربع الثاني لتحصل على درجة عالية عن الحرف C ولكنك ستنال عقوبة بالخصم للانتقال ولكن ستجد أن الحرف التالي في الترتيب رقم ١٢ هو T فتنتقل مرة ثانية للمربع الأول للحصول على الدرجة العالية لكنك ستنال عقوبة الانتقال مرة ثانية، وعند بلوغ الحرف رقم ١٥ تنتقل للمربع الثاني حيث معظم الأحرف من رقم ١٥ إلى رقم ٢٥ من النوع C أو G، تقابلك نفس الحالة عند الحرف رقم ٢٠ حيث أنه A (يمكن تسمية هذه الطريقة بالإستراتيجية الأولى). أما الإستراتيجية الثانية هي أن تبقى في المربع الأول من الحرف رقم ١ حتى رقم ١٥ دون تحرك متحملا الدرجة المنخفضة للحرف رقم ١١ لكن متجنباً عقوبة الخصم مرتين، ومن الحرف رقم ١٥ حتى النهاية فتنتقل للمربع الثاني متحملا الدرجة المنخفضة عند الحرف رقم ٢٠ ولكن متجنباً الخصومات. ويتضح أن الإستراتيجية الثانية هي التي ستعطيك أعلى الدرجات لكسب الدور.

فيما سبق حاولنا تبسيط الأمر بقدر المستطاع لتفهم المهام المطلوبة من برامج ماركوف الإحصائية، لكن مع الأخذ في الاعتبار أن رقعة المربعات ستكون من عدد كبير من المربعات وعليها أن تختار المسار path الأمثل بينها للوصول للنتيجة المرجوة، وتخيل حجم التباديل والمحاولات الواجب القيام بها لإتمام الدور. ومما زاد من الطين بلة أن

برامج ماركوف هذه يجب أن يكون بها قدر من المعلومات الخفية على اللاعب، لذلك تسمى برامج ماركوف الخفية Hidden Markov Models وتعرف اختصاراً HMM. ومثال ذلك أنك ستبدأ اللعب دون أن تعرف مكان الأبتداء أو مكان الإنتهاء، كمن يلعب أحد ألعاب التسلية (السلم والثعبان مثلاً) في الظلام وأن شخصاً آخر يلعب نيابة عنه دون أن يطلعه على ما يجري من أحداث وعليه جمع النقاط فقط !!!

١.٢.٤. برامج ماركوف و الوراثة Markov Models & Genetics

لدراسة تتابع ما من الـ DNA باستعمال برامج ماركوف، يجب أن نحدد رموز (حروف) الموضوع تحت الدراسة، وهو في هذه الحالة معروفة لدينا، فهي النيوكلووتيدات الأربع (A,G,C,T) التي يمكن أن تتواجد أو تتبادل مواقعها على هذا التتابع، وشكل (٥-٤) يوضح تلك الرموز والتبادلات الممكنة لكل منها (قدرها ١٦ لكل منها)، ومع تجاهل البداية (start) والنهاية (end) في هذه المرحلة.



شكل (٥-٤) : رسم توضيحي للنيوكلووتيدات الأربع واحتمالات التبادل والإحلال بينها.

ومن معلوماتنا الوراثية نعلم أن الجينومات في الكائنات الراقية تتميز بوجود جزر من تتابعات المتكررة من C_m G islands (C_m G) ولكن نعلم أيضاً أن الجينوم يتعرض لعمليات الـ methylation الدائمة التي تحول السيتوسين إلى الثيمين أي أن مثل تلك الجزر ستصبح C_m G islands، كما نعلم أن مناطق البروموتورات (promoters) لا تحدث فيها عمليات الـ methylation، أي تبقى بها الجزر من C_m G دون تغير. وحيث

أن احتمال كل رمز في موقع ما S_p سيتوقف فقط على احتمال الرمز الذي قبله في الموقع $p-1$ أى S_{p-1} ، وباستعمال تتابعات معروف أنها غنية في جزر G C_p (overabundant model) و تتابعات أخرى معروف أنها فقيرة في تلك الجزر (underabundant model)، يمكننا حساب الاحتمالات لتبادل الرموز بينها البعض وذلك بناء على المشاهدات التجريبية، والجدول التالي يلخص تلك الاحتمالات من نتائج النموذج الفقير في تلك الجزر:

From / To	A	C	G	T
A	0.300	0.205	0.285	0.210
C	0.322	0.298	0.078	0.302
G	0.248	0.246	0.298	0.208
T	0.177	0.239	0.292	0.292

والجدول التالي يلخص الاحتمالات المتحصل عليها من نتائج النموذج الغنى بتلك الجزر:

From / To	A	C	G	T
A	0.180	0.274	0.426	0.120
C	0.171	0.368	0.274	0.188
G	0.161	0.339	0.375	0.125
T	0.079	0.355	0.384	0.182

مما سبق من بيانات، أصبح لدى البرنامج نموذجين للقياس، الأول نموذج يمثل التتابعات الفقيرة في تلك الجزر (underabundant model)، والثاني نموذج يمثل التتابعات الغنية في هذه الجزر (overabundant model). وحيث أن احتمال حدوث رمز S عند الموقع P (S_p) يتوقف فقط على الرمز الذي قبله عند الموقع $P-1$ (S_{p-1}) وليس على جملة التتابعات السابقة، ولحساب احتمال أن تتابع ما يتوافق (fits) مع نموذج ما يقرر حاصل ضرب الاحتمالات الشرطية (conditional probabilities) كالتالي :

$$.P(x) = P(x_L/x_{L-1}) P(x_{L-1}/x_{L-2}) \dots P(x_2/x_1) P(x_1)$$

والتي تمثل رياضيا بالمعادلة التالية:

$$P(x) = P(x_1) \prod_{i=2}^L a_{x_{i-1}x_i}$$

حيث $a_{x_{i-1}x_i}$ هو قيمة الاحتمال لانتقال الرمز من الموقع $i-1$ إلى الموقع i .

ولسهولة الحسابات دعنا نقول بأن الاحتمالات المحسوبة من النموذج الفقير في $C_p G$ - [underabundant model]

$$P(A) = P(T) = 0.3 \quad \& \quad P(C) = P(G) = 0.2$$

وبأن الاحتمالات المحسوبة من النموذج الغني في $C_p G$ [overabundant model] =

$$P(A) = P(C) = P(G) = P(T) = 0.25$$

الآن دعنا نحاول دراسة التتابع **GGCGACG** على سبيل المثال لتحديد مع أي من النموذجين يتوافق. واحتمال هذا التتابع تبعاً للـ [underabundant model] هو،

$$P(G)P(G/G)P(C/G)P(G/C)P(A/G)P(C/A)P(G/C)$$

والذي يساوي-

$$(0.20)(0.298)(0.246)(0.078)(0.248)(0.205)(0.078) \\ = 0.000000453499$$

واحتمال هذا التتابع تبعاً للـ [overabundant model] يساوي -

$$(0.25)(0.375)(0.339)(0.274)(0.161)(0.274)(0.161)(0.274)(0.274)(0.125) \\ = 0.0010526$$

وبناء على تلك الحسابات - يمكن القول بأن هذا التتابع أكثر أرححية في أن يكون منتمى للنموذج الثاني أي غني في $C_p G$ ويتبع الـ overabundant model.

ويتضح لنا من حسابات الاحتمالات السابقة، أن القيم تؤول للصفر (zero) بسرعة كبيرة وللتغلب على ذلك يستعمل التحويل اللوغارتمى أو ما يسمى بـ log statistics.

فى استخدامنا السابق لطريقة ماركوف كنا نبحث عن التوافق أو الإنتماء، ولكن يمكن أن نستخدم هذه الطريقة أيضا للتفرقة (discrimination) بين المعطيات. فمثلا إذا ما تساءلنا عن مدى الاختلاف بين النموذجين السابقين أى الـ overabundant والـ underabundant. وإذا ما كانا غير مختلفين بالقدر الكافى عن بعضهما البعض، فلن تتوفر لدينا معلومات حقيقية وكافية لتقدير أرجحية أن تنابع ما فى أن ينتمى لى منهما. ولإختبار قدر الاختلاف بينهما نستخرج قيم اللوغارتم للأساس الثانى لنسب الاحتمالات للنموذجين من الجدولين السابقين وتحسب القيم كما هو فى الجدول التالى :

From / To	A	C	G	T
A	-0.740	0.419	0.580	-0.803
C	-0.913	0.302	1.812	-0.685
G	-0.624	0.461	0.331	-0.730
T	-1.169	0.573	0.393	-0.679

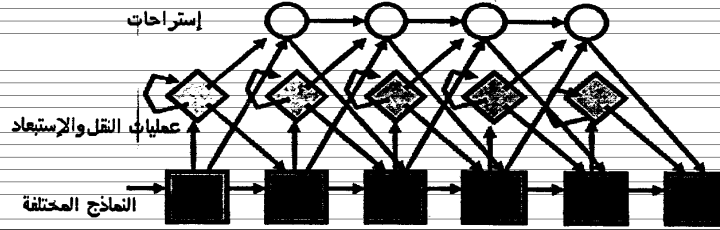
ومن القيم السابقة يتضح أن الفرق واضح بين كلا النموذجين فالقيم التى ستتنمى للنموذج underabundant model ستكون سالبة القيمة أما القيم التى ستتنمى الى النموذج overabundant model ستكون موجبة، مما يؤكد أنهما نموذجين مختلفين.

٢.٣.٤ نماذج ماركوف الخفية Hidden Markov Models

نماذج ماركوف التى تناولناها بالشرح والتفسير فيما سبق محدودة الإستعمال، وعادة تكون متطلبات البحث والتدقيق فى المعلومات الوراثية هتداول أكثر من نموذج أو ما يسمى بالوضعىة (state)، فلمعرفة أرجحية إنتماء بروتين ما لأحد العائلات البروتينية (protein families) نحتاج أكثر من نموذج قد يصل عددها إلى عشرة، كذلك تكون التتابعات تحت الدراسة أكبر بكثير من تلك التى إستعملناها فى أمثلتنا السابقة. وفى حالة النموذجين (الغنى و الفقير بالجزر من G و C)، لكى نبحث عنهما فى مقطع كبير من الـ DNA الجينومى، فلن نعرف أماكن تلك الجزر على وجه التحديد (أى خفية علينا)،

لذلك وجب دمج النموذجين مع بعضهما مع السماح للبرنامج بالتنقل خلال مسار محدد (path) وعادة يكون خفي لحساب احتمالات التحول (transition) واحتمالات الإستبعاد (emission) لكل حدث أو رمز.

والتخطيط المبين بشكل (٦-٤) يبين تصميم لبرنامج إحصائي من نوع ماركوف الخفي (HMM)، حيث يتم تحديد النماذج القياسية المختلفة التي يمكن قياس التتابع تحت الدراسة عليها، وبرنامج تنفيذ عمليات المسح من تحولات (transitions) وإستبعادات (emission)، أما مناطق الأسرحة بين التحولات والإستبعادات لمسار (path) ما من ضمن المسارات المحتملة فيجب تواجدها كمناطق محايدة لسهولة تحديد المسار الأمثل.



شكل (٦-٤) : رسم تخطيطي لتصميم أحد برامج ماركوف الخفية (HMM).

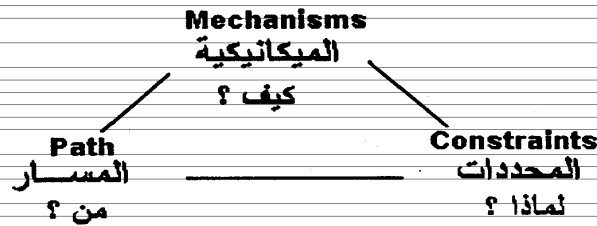
ويتوفر على شبكة المعلومات العالمية INTERNET وعلى الموقع الخاص بمدرسة الطب بجامعة واشنطن (Washington University in St. Louis, School of medicine) عدة برامج لإستعمال برامج ماركوف الخفية للدراسة تتابعات البروتينات والمعروفة بأسم HMMER وآخر نسخة منها صدرت سنة ٢٠٠٤ باسم hmmer 2.3.2

الموقع: <http://hmmer.wustl.edu>

ومع عظمة هذه الأساليب الحسابية الحديثة وإسهاماتها الجليلة في تفهم مواضيع مثل البيولوجيا الجزيئية والجينومكس إلا أن المبدأ الفيزيائي العام، ألا وهو "عدم التأكد - Uncertainty principle" يمكن أن يلقي بظلاله عليها، بالاعتقاد بأن

محاولة التنبؤ بالسلوك البيولوجي أو الوراثة في الطبيعة هي عملية غير مؤكدة بل
مستحيلة !!

دعنا الآن نتناول موضوعا وراثيا مرتبط بموضوع دراستنا وقد يقع تحت طائفة
مبدأ عدم التأكد هذا، ألا وهو دراسة التطور بين الكائنات اعتمادا على دراسة مدى التشابه
والإختلاف بين البروتينات لكائنات مختلفة تربطها درجات من القرابة، أو ما يسمى
بـ **فيولوجينيا** البروتين (protein phylogeny). ان دراسة تتابعات الأحماض الأمينية
للسلاسل الببتيدية في الكائنات المختلفة لإستخلاص العلاقات الفيولوجينية منها يعتمد
على ثلاثة محاور رئيسية متداخلة : أولاها هو **ميكانيكيات** (mechanisms) التغيير
لدراسة أسباب منشأ التغيرات بين البروتينات (تحاول الأجابة على كيف؟ حدث التغيير) و
ثانيها هو **المسار** (path) الذي سلكته، لدراسة مسار التغيير بين التشابه والمتباعد (تحاول
الأجابة على ممن؟ جاءت الإختلافات) وثالثها هو **المحددات** (constraints) لهذه التغيرات
لدراسة العوامل التي قد تساعد أو تعاكس ظهور هذه الإختلافات (تحاول الأجابة على لماذا؟
توجد الإختلافات). والرسم التوضيحي التالي يبين هذه المحاور الثلاث.



فكل محور منها يؤثر على الآخر بدرجة أو بأخرى، ولأكمال الدراسة ولزيادة
مصداقيتها لابد من تناولها جميعا في الحساب عند دراسة فيولوجينية بروتين ما ضمن
عدد من الكائنات، ولكن القصور المادي لتصميم التجارب يحول دائما دون هذا. فعند دراسة
الفيولوجيني أو المسار فلا بد أن نفترض نمودجا ثابتا للمحورين الآخرين، مثل إفتراض
ميكانيكية محددة (من ضمن احتمالات عديدة) وكذلك يجب إفتراض وجود نمودجا
للإنتخاب محدد (من ضمن احتمالات عديدة). بناء على هذا التبسيط التجريبي
(القصور) فإن النتائج التي نتوصل إليها عادة غير مؤكدة وبالتالي تندرج تحت طائفة مبدأ

عدم التأكد لـ Heisenberg. لتوضيح هذه النقطة دعنا نقارن بين تنبؤات الأحماض الأمينية لثلاثيات الببتيدات الافتراضية التالية من أربعة كائنات مختلفة:

. . VGM . .

. . VAM . .

. . VPM . .

. . VLM . .

فالإختلافات عند الموقع الأوسط يمكن تفسيرها تبعاً لواحد من الاحتمالات الثلاث التالية:

- قد يكون الكودون الخاص بالحمض الأوسط يمثل نقطة ساخنة (hot spot) للطفور. هنا ركزنا على "الميكانيكية" دون المسار والمحددات.
- قد يمثل ذلك خطوات المنشأ من سلف مشترك بالترتيب التالي $G > A > P > L$ حيث أن هذا الاحتمال ينتج من إستبدال نيوكلويدة واحدة بينما باقى الاحتمالات تنشأ من إستبدال أكثر من نيوكلويدة واحدة. هنا ركزنا على "المسار" دون الميكانيكية والمحددات.
- قد يمثل ذلك عدم وجود محددات (مثل الإنتخاب الطبيعي) عند الكودون الأوسط. هنا ركزنا على "المحددات" دون الميكانيكية والمسار.

٤.٤.٤. ملحق (1) Appendix.

البرمجة (computer programming) تعتبر الآن جزءاً أساسياً من المعلوماتية الحيوية، وإن كانت ليست من مهام البيولوجيين وتعتمد على المتخصصين في علوم الكمبيوتر، إلا أن هذا الملحق يتناول بشكل سريع بعض من المعلومات الأساسية لأحد لغات البرمجة للكمبيوتر، وهى لغة Perl، الأكثر استعمالاً فى البرمجة الوراثية وعسى أن تكون النواة ليزداد اهتمام أحد البيولوجيين بالموضوع ثم التخصص فى هذا المجال الجديد.

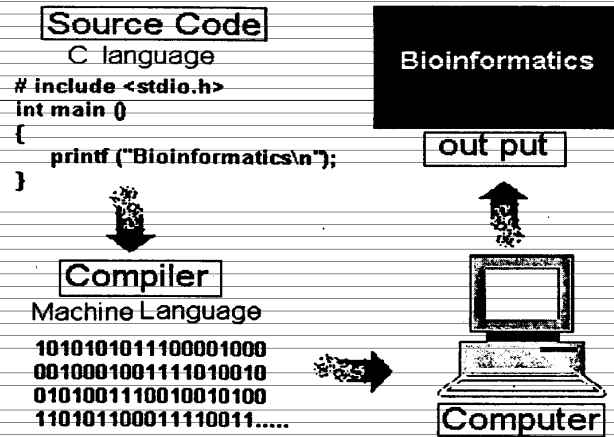
لغات البرمجة للكمبيوتر عديدة وفي كل يوم هناك الجديد في هذا المجال، وعلى تعدد واختلاف تلك اللغات فلكي يتفهمها الكمبيوتر فيجب أن تحول جميعها إلى لغة واحدة خاصة بالكمبيوتر وتعرف باسم بلغة الآلة (machine language) وهي لغة ثنائية (binary) كما هو موضح بشكل (٧-٤). ولغات البرمجة أو ما يعرف بأسم source codes متعددة الأنواع ويمكن إجمالها في الآتي،

١. لغات المعالجة (compiled languages) – مثل لغات C, C++, Java وهي تمتاز بسرعة العمل والتنفيذ ولكنها صعبة في الكتابة وتستخدم في برامج أنظمة التشغيل (operating systems) وبرامج معالجة الكلمات (word processing).

٢. لغات النصوص (scripting language) – مثل لغات Perl, Python, Ruby, Tel وهي تمتاز بسهولة الكتابة ولكن بطيئة في العمل والتنفيذ لذلك فهي مناسبة للبرامج الصغيرة.

٣. لغات الشبكة العالمية للمعلومات (Web languages) – مثل html, PHP, Java script وهي خاصة بالربط بين صفحات الشبكة المتداخلة وربطها مع قواعد المعلومات.

٤. لغة قواعد المعلومات (databases language) – مثل SQL - لغة قادرة على تداول المعلومات بين قواعد المعلومات المختلفة وشبكة المعلومات العالمية



شكل (٧-٤) : تخطيط يمثل خطوات تداول برنامج مكتوب بلغة C وتحويله للغة ثنائية binary حتى يفهمه الكمبيوتر وينفذه.

واللغة الأكثر استعمالاً في مجال المعلوماتية الحيوية هي لغة Perl وهي اختصار لـ Practical Extraction and Report Language وقد اقترحها Larry Wall لأنها مناسبة لتداول النصوص واستخلاصها من شبكة المعلومات العالمية والقدرة على تلقي وإرسال المعلومات بسهولة وأن كان يعيبها بطء التنفيذ. وفيما يلي بعض من القواعد الأساسية لهذه اللغة:

قواعد بناء الجمل Syntax rules:

1- Statements are terminated by a semi-colon

- Print ("Hello!\n");

2- Text blocks are determined by curly brackets

- If (\$a == \$b) {
print ("a=b!\n");

}

3- Comments are indicated by sharp sign

`$a = 10; # set $a equal to 10`

4- Separate variable names with spaces, otherwise space has no meaning

- `$a + $b`; is the same as `$a + $b`;

5- Common conventions: `\n` = new line, `\t` = tab, `" "` = string

6- **Order matters**: statements are evaluated in descending order.

المهام Assignments

1- Equal sign represent variable assignment

- `$A = B`

2- Binary assignments operators:

- `$A = $A + 5; => $A += 5;`

- `$B = $B - 6; => $B -= 6;`

انواع المتغيرات Variable types

1- Dollar-sign (\$) variable represents a scalar (number) or string

- `$DNA_length = 10;`

- `$DNA_sequence = "ATTAGCCGAT";`

2- At-sign (@) variable represents an array, (\$) sign represents individual array element

- `@DNA = (A,T,T,A,G,C,C,G,A,T);`

- `$DNA[0]` is equal to "A"; `$DNA[1]` is equal to "T";

3- Percent-sign (%) represents a hash, (\$) represent individual hash element

- `%DNA = ("First" => "A", "Second" => "T");`

- `$DNA["First"]` is equal to "A";

، Arithmetic & Logical operations # العمليات الحسابية والمنطقية

- 1- Addition: \$a = 5 + 6; # \$a equal 11
- 2- Subtraction: \$a = 6 - 5; # \$a equal 1
- 3- Multiplication: \$a = 3 * 2; # \$a equal 6
- 4- Division: \$a = 6 / 2; # \$a equal 3
- 5- Modules: \$a = 3 % 2 = 1
- 6- Auto-increment: \$a = 0; \$a++; # \$a is equal 1
- 7- Identify test: if (\$a == 6); # \$a is equal to 6
- 8- Not equal test: if (\$a != 6); # \$a is not equal to 6
- 9- Less than, greater than: \$a < 6; \$b > 5
- 10- AND operator: \$a == 6 && \$b > 4; # \$a equal to 6 AND \$b is greater than 4
- 11- OR operator: \$a == 6 || \$b < 4; # \$a equal to 6 OR \$b is less than 4

، Conditional statements # الأوامر الشرطية

- 1- if/else statements:
 - if statement {do if statement is true}
 - else statement {do if statement is false}

: Strings # الأوتار

- 1- String connection is a period (.)
 - \$a = "Hello"; \$b = "World";
\$c = \$a. \$b; # \$c is "Hello World"
- 2- String length: \$length = length(\$string)
 - \$text_length = length(\$c); # \$text_length is 10
- 3- String reverse: \$rev_string = reverse(\$string);
 - \$rev_c = reverse(\$string); # \$rev_c is "dlro WolleH"

المتعلقات Loops

1- for statements:

```
- for($i = 0; $i < 20; ++){
    $DNA[$i] = "A";
}
```

2- foreach statements:

```
- foreach $i (@DNA) {
    $i = "A";
}
```

3- while statement:

```
- while ($i < 20) {
    $DNA[$i] = "A";
    $i++;
}
```

المعطيات والمخرجات Input/Output

1- Standard input: <STDIN>

2- Standard output: print

3- Opening file: open (FILEHANDLE,"filename");

```
- open (DNASEQ, "dnaseq.txt file");
```

4- Reading file: single line => \$DNA = <FILEHANDLE>;

```
whole file => @DNA = <FILEHANDLE>;
```

المصفوفات Array

1- Split function: splits a string into an array of letters.

```
$seq = "ATAGCCAT");
```

```
@DNA = split(/,$seq); # $DNA[0] is "A", $DNA[1] is "T"
```

- 2- **Push/Pop:** Push adds value to end of array, pop removes last value of array.

```
push (@DNA, "G"); # @DNA is {A,T,A,G,C,C,A,T,G}
```

```
$last = pop (@DNA); # $last is "G"
```

- 3- **Reverse:** reverses order of the array

```
@DNA = reverse (@DNA); # @DNA is {T,A,C,C,G,A,T,A}
```

- 4- **Length of array:** scalar @array

```
$size of array = scalar @DNA; # $size of array is 8.
```

تعبيرات عامة Regular Expressions

- 1- **Matching:** \$string =R/pattern/

```
$DNA = "ATATAAAGA";
```

```
if ($DNA =R/TATA/ {
```

```
    print ("Contains TATA element\n");
```

```
}
```

- 2- **Substitution:** \$string =Rs/pattern/replacement pattern/(g);

- 3- **\$DNA =Rs/TATA/GGGG/g; # \$DNA is now**

```
"AGGGGAAGA"
```

- 4- **Wildcards:** [ATGC] matches A or T or G or C;

```
[^0-9] matches all non-digit characters;
```

```
A{1,5} matches a stretch of 1 to 5 "A" characters
```

فيما سبق حاولنا التعرف على أهم القواعد للغة Perl ولزيت من التفصيل يمكن

اللجوء إلى المواقع المتخصصة في هذه اللغة كالتالية:

www.perl.com/perl

www.perl.com/cpan

مثال: البرنامج الأساسي لدراسة توافق تتابع ما بلغة Perl.

```
for (i = 0 to l = m - n) {
    j = 0;
    while (j < n and Qj == Ti + j) {
        if (j == n - 1) {
            return l;
        }
        j = j + 1;
    }
}
```

where: Q = DNA sequence (Query)

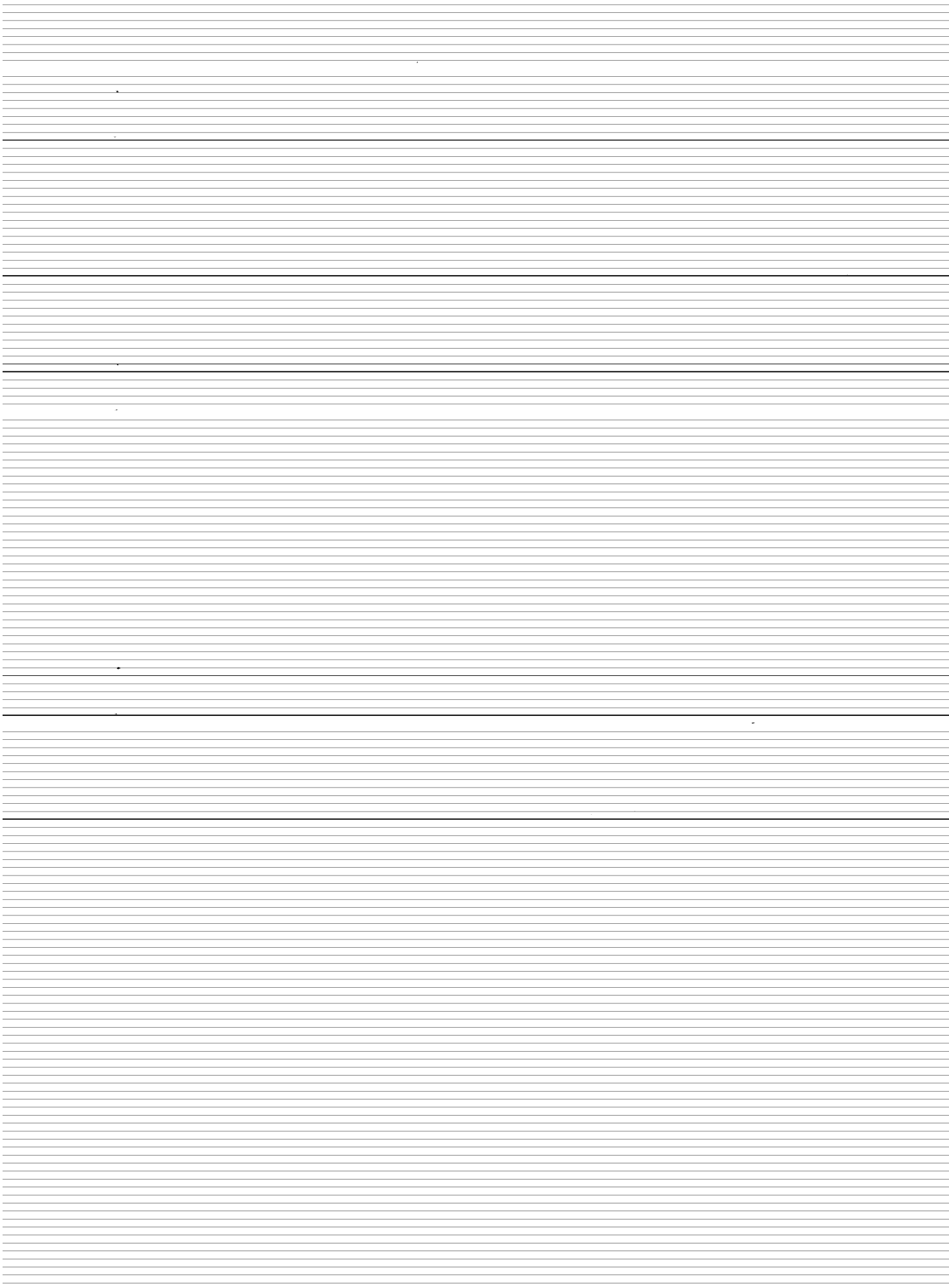
T = Chromosome sequence (Target)

Q_j = Nucleotide in Query sequence at position j

T_i = Nucleotide in Target sequence at position i

n = length of Q

m = length of T



٥. قواعد المعلومات

البيولوجية

Biological Databases

إعداد : آمال عبد العزيز

تعتمد المعلوماتية الحيوية اعتمادا جوهريا على توفر قواعد المعلومات البيولوجية بمفهومها الحديث. وقاعدة المعلومات في أبسط صورها قد تكون ملفا صغير يحفظ عدة بيانات عن موضوعا ما أو قد تكون عبارة عن مخزن كبير (أرشيف) لحفظ ملفات عن قطاع من المواطنين و الذى قد يشغل عدة أدوار من مبنى كبير، ومع ان هذه المحفوظات هي صورة من صور قواعد المعلومات، إلا أنها ليست المقصودة في مجال اهتمامنا الحال. فقواعد المعلومات (databases) التي نقصده هنا هي : مجموعة البيانات أو المعلومات المخزنة والمبوية بطرق إلكترونية باستعمال تقنيات الكمبيوتر وتتميز بمدخل أمامي (front-end) حيث يمكن للمستعمل الدخول والبحث والتدقيق بالقاعدة، ويجب أن يتوفر لها أيضا مدخل خلفي (back-end) حيث يمكن للمشغل أو المسئول عن القاعدة أن يتناول بيانات القاعدة بالإضافة و الحذف أو التحديث (up-dating). ويرجع الفضل إلى مهندس الكمبيوتر الشهير Charles Backman في تصميم أول تلك القواعد بمفهومها الحديث مع مطلع الستينات من القرن الماضي. ويمكن اعتبار الأطلس الذى نشرته العائلة Margaret Dayhoff وزملائها سنة ١٩٦٥ بعنوان "Atlas of protein sequences and structure"، والذي كان النواة التي أصبحت فيما بعد قاعدة المعلومات عن البروتينات والمعروفة بأسم PIR database وهو قد يعتبر أول تلك المحاولات الحديثة.

وقواعد المعلومات البيولوجية بالنسبة للمهتمين بالمعلوماتية الحيوية توفر لهم المعلومات الأساسية والتي بدونها تصبح الدراسة مستحيلة، ويمكن إجمال فوائد قواعد المعلومات البيولوجية في الآتي :

(١) توفير المعلومات الجينية والبروتينية للباحثين.

(٢) تخزين وتبويب المعلومات الجينية والبروتينية على هيئة صفحات مقروءة بالكمبيوتر.

٥. ١. طرز قواعد المعلومات البيولوجية - Types of biological databases

- عادة تقسم قواعد المعلومات البيولوجية لطرز عديدة تتبع فيها معايير مختلفة، وفيما يلي أهم الطرز حسب تلك المعايير:
- (١) بناء على نوع الكائن - فقد تقسم القاعدة بناء على تخصصها عن معلومات لكائن واحد فقط، فقد تختص قاعدة بجينات الإنسان وأخرى بجينات الفأر وثالثة بجينات نبات الأرز وهكذا.
 - (٢) بناء على نوع البيانات - وهناك قواعد تتميز بنوع البيانات التي تختص بها، فهناك قواعد لتتابعات النيوكلووتيدات و أخرى لتتابعات البروتينات وثالثة لبيانات الأبعاد الثلاثية (3D) للجزيئات ورابعة للتعبير الجيني وهكذا.
 - (٣) بناء على طبيعة البيانات - هل البيانات في القاعدة تمثل النتائج الأولية (القواعد الأولية) أم تمثل النتائج بعد التحليل والتدقيق (القواعد الثانوية).
 - (٤) بناء على الإتاحة (Availability) - تقسم القواعد أيضا على مدى إتاحتها على الشبكة الدولية للمعلومات، فمنها ما هو متاح للجمهور مجانا ومنها ما هو متاح مجانا مع الحفاظ على حقوق الملكية الفكرية ومنها ما هو متاح للمتخصصين أو المشتركين فقط وهكذا.
 - (٥) بناء على التصميم - يعتبر التصميم الفني للقاعدة هو أهم المعايير التي تميز قاعدة عن أخرى، وفيما يلي نبذة مبسطة عن نظم تصميم وبرمجة قواعد المعلومات.

١.١.٥. نماذج تصميم قواعد المعلومات Modeling Databases.

يعتبر تصميم وبرمجة قواعد المعلومات من المهام الأساسية لإنشاء تلك القواعد، وهناك العديد من الطرق وبرامج الكمبيوتر المتاحة لهذا الغرض، ولكن أهم ميزتين يجب أن يتصف بهما أي نموذج لتصميم قاعدة ما، هما،

- البساطة Simplicity.
- تجنب تكرار البيانات Redundancy.

وفي مبدأ الأمر كانت النماذج بدائية لا تختلف كثيرا عن الجداول أو ما يسمى spreadsheet ويطلق على هذه القواعد النماذج المسطحة (ذات البعدين) flat-files. والجداول التالية يوضح أحد تلك النماذج لتسجيل بيانات الطلبة.

الاسم	الرقم الكودي	المقرر	التقدير
محمد محمد محمد	٤٤٤٤٤	ورائة (١٠١)	ممتاز
محمد محمد محمد	٤٤٤٤٤	ورائة (١٠٢)	جيد جدا
محمد محمد محمد	٤٤٤٤٤	ورائة (١٠٣)	جيد جدا
سامي سامي سامي	٦٦٦٦٦	ورائة (١٠١)	جيد جدا
سامي سامي سامي	٦٦٦٦٦	ورائة (١٠٢)	جيد
سامي سامي سامي	٦٦٦٦٦	ورائة (١٠٣)	ممتاز

ويلاحظ أن مثل هذه القواعد به كمية كبيرة من البيانات المتكررة والبحث فيها يستغرق وقتا كبيرا، لذلك لجأ المصممون لتجزئة البيانات وتقليل التكرار عن طريق تصاميم سميت بالنماذج الهرمية أو الشجرية Hierarchical (tree) model حيث ترتب جداول المعلومات في أنظمة هرمية مرتبطة ببعضها البعض كما هو موضح بالجداول التالية:

الجدول الأول

الاسم	الرقم الكودى
محمد محمد محمد	123-44444
تلى تلى تلى	143-66666

الجدول الثانى

رقم المقرر	اسم المقرر	الوحدات	الفصل	القسم
١٠١	أساسيات الوراثية	٣	١	الوراثية
١٠٢	علم الأحياء	٤	٢	الوراثية
٥	صيرولوجى النبات	٤	١	النبات
١٠٣	تقسيم حيوان	٣	١	للمحورين
١١١	أعمال نووية	٣	١	للكيمياء

الجدول الثالث

القسم	رئيس القسم	عدد المحفلات المتاحة
الوراثية	أ.د. على على	٢٣
النبات	أ.د. هدى كامل	٣٣
للمحورين	أ.د. حسن حسن	٢١
للكيمياء	أ.د. محمد محمد	٤١

وهكذا.....

حيث يقل تكرار البيانات بالجدول وتوفر معلومات أكبر ولكن مثل هذه التصميم قل استعمالها الآن، وتصمم القواعد الآن باستخدام برامج الكمبيوتر اعتماداً على الأسس الرياضية والجبرية مثل نماذج Relational model التى تعتمد على لغة SQL وهى اختصار Structured Query Language وهى لغة من لغات الكمبيوتر توفر طرق البحث والتدقيق والتحديث لقواعد المعلومات. وحديثاً تستعمل لغات أخرى لإنشاء قواعد المعلومات خصوصاً التى تتداول النصوص والوسائل المرئية والسمعية عبر شبكة المعلومات العالمية internet مثل لغات FTP, HTML, CORBA, XML وغيرها.

٢.٥. قواعد المعلومات الوراثية Genetic Databases.

يتواجد الآن عدد كبير من قواعد المعلومات المتخصصة لتبويب وتجميع البيانات الخاصة بالأحماض النووية والبروتينات، وهناك عدد محدود منها يمثل أهم تلك القواعد وأكثرها شهرة واستعمالاً. ففىما يخص تتابعات الأحماض النووية فإن القاعدة التابعة لـ NCBI الأمريكية و قاعدة EMBL الأوروبية وقاعدة DDBJ اليابانية هم الأهم والأكثر استعمالاً، وفيما يخص البروتينات فإن قاعدة SWISS-PROT الأوروبية هى الأهم. بالإضافة لتلك القواعد فهناك مئات من القواعد الأخرى الأقل أهمية ولكنها قد تكون أكثر تخصصاً، كل فى مجاله، والجدول التالى يلخص بعض من أهم تلك القواعد:

القاعدة	الموقع	ملاحظات
VectorDB	http://vectordb.atcg.com/	تحتوى على العديد من تتابعات أدوات النقل
Codon usage	http://www.dna.affrc.go.jp/	أنواع الكودونات في الكائنات المختلفة.
TRANSFAC	http://transfac.gbf.de/	تتابعات عدد من <i>cis</i> -acting DNA تعمل كمُنظمات للنسخ.
RNA world	http://www.imb-jena.de/	تحتوى على وصلات لقواعد أخرى بها عديدة من تتابعات الـ RNAs الصغرى.
NDB	http://ndb-mirror-2.rutgers.edu/	تجميع البيانات عن الأحماض النووية.
PRINTS	http://www.biochem.ucl.ac.uk/	تشمل مجموعة من التراكيب البروتينية المحفوظة، لتحريف البروتينات.
Pfam	http://www.sanger.ac.uk/	تشمل على أجزاء بروتينية .domains
GSDB	http://www.ncgr.org/gsdb/	بها قدر ممتاز من التتابعات الجينومية.
REBASE	http://rebase.net.com/	تشمل معلومات عن الإنزيمات المحددة ومناطق تعارفها.
GCRDb	http://www.gcrdb.uthscsa.edu/	تحتوى معلومات وافية عن الـ G-proteins.

١.٢.٥. المركز القومى الأمريكى للمعلومات البيوتكنولوجية NCBI .

تنبّهت الحكومة والكونجرس في الولايات المتحدة الأمريكية في نوفمبر سنة ١٩٨٨ إلى الحاجة الملحة لإستخدام الكمبيوتر لتخزين وتبويب المعلومات المتنامية في مجالات الكيمياء الحيوية البحتة و الكيمياء الحيوية الطبية، لذلك عمل على سن القوانين التي تساعد في عمل مكتبة قومية خاصة بالكيمياء الحيوية الطبية وأن تكون تابعة للمكتبة الطبية القومية NLM التابعة بدورها للمعهد القومى للصحة NIH، وذلك للعمل على تحسين وتوفير الأدوات التحليلية الجديدة وللمساعدة في إستيعاب المفاهيم الجديدة في الوراثة الجزيئية و بالتالى معرفة دورها في دراسة الحالات المرضية. وعليه أنشئ المركز القومى للمعلومات البيوتكنولوجية NCBI حيث حدد له أربعة مهام رئيسية وهي :

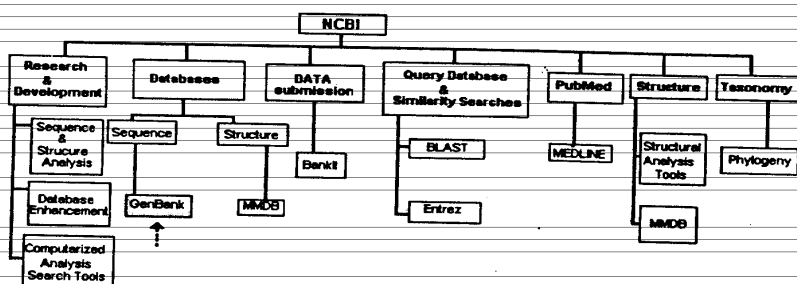
١. ابتكار أدوات أوتوماتيكية تستطيع تحليل و تخزين المعلومات الخاصة بالبيولوجيا الجزيئية عموما و الوراثة الجزيئية و الكيمياء الحيوية خصوصا.

٢. تسهيل استخدام البرامج التحليلية المتاحة للإتصال بالبيانات الرئيسية أو الأساسية (على سبيل المثال إتاحة تبادل المعلومات بين العلماء المهتمين بالنواحي الطبية).
٣. العمل على ربط الجهود العالمية مع بعضها البعض في مجال البيانات البيولوجية.
٤. العمل على تسهيل الاتصالات البحثية لتوفير وسائل التحليلات التركيبية و الوظيفية للبيولوجيا الجزيئية.

وقد بدأ الـ NCBI متواضعا ولكن مع مرور الأيام أصبح أكبر قاعدة للمعلومات في العلوم البيولوجية الجزيئية، فهو يقدم الآن خدمات كثيرة من خلال مواقع متخصصة على الشبكة العالمية للمعلومات (الأنترنت)، ويمكن إجمالها في الخدمات التالية:

1. PubMed (Public MEDLINE).
2. BLAST (Basic Local Alignment Search Tool).
3. Entrez.
4. BankIt (World Wide Web Submission).
5. OMIM (Online Mendelian Inheritance in Man).
6. Taxonomy.
7. Structure.

وشكل (١-٥) يلخص أهم تلك المهام والخدمات.



شكل (١-٥) : رسم توضيحي يمثل المهام والخدمات التي يوفرها المركز القومي للمعلومات البيوتكنولوجية NCBI.

فمن خلال الموقع Entrez أو الموقع BLAST أو حتى مباشرة من موقع NCBI يمكن الدخول إلى الموقع المعروف باسم "بنك الجينات" GenBank هو بدون شك الأهم لنا نحن الوراثةيون.

١.١.٢.٥. موقع بنك الجينات GenBank.

يتضمن هذا الموقع تتابعات النيوكليوتيدات لكل من الـ DNA والـ RNA، حيث تحفظ البيانات المختصة والتي يتم تجميعها من الباحثين مباشرة أو من المعامل العلمية المختصة أو من مكاتب البراءات العلمية Patent Offices وكذلك يتم ضم المعلومات من القواعد الأوروبية EMBL واليابانية DDBJ على أسس يومية كجزء من التعاون الدولى فى هذا المجال. وقد بلغت المعلومات المخزنة فى هذا الموقع فى أغسطس ٢٠٠٤ حوالى ٤١,٨ مليون قاعدة من حجم تتابعات قدره ٢٧,٢ مليون تتابع تمثل جينومات العديد من الأنواع حقيقية وغير حقيقية النواة. وتقدر الزيادة التى تمت على البنك خلال الفترة المنصرمة من سنة ٢٠٠٥ بحوالى ٧,٩ مليون قاعدة جديدة كما ذكر Benson et al. سنة ٢٠٠٥.

يمكننا دخول هذا البنك الجينى عن طريق الموقع Entrez حيث يوجد الموقع GenBank والذي بدوره تبويب فيه البيانات تحت العديد من تحت الأقسام : ومن أهمها تحت قسم التقسيم sequence-based taxonomy أى حسب نوع وأسم الكائن العلمى، حيث يحتوى الآن على حوالى ١٦٥,٠٠٠ نوع مختلف وتقدر الزيادة بهذا القسم بحوالى ٢٠٠٠ نوع كل شهر - أو من خلال تحت قسم مصدر التتابع حيث تتراكم يوميا البيانات الجديدة من مختلف التتابعات لكن تبقى التتابعات المستخلصة من EST هى الأكثر حتى الآن فهى تمثل ٢٩٪ من حجم البيانات، ولكن يمكن منها الوصول إلى بيانات الجينات المحددة UniGene والتي تبلغ حوالى ٧٠٠,٠٠٠ جين مستقل يمثلون ٥٠ كائن مختلف.

عند طلب أى معلومات من هذا الموقع GenBank عن تسلسل تتابعات جين أو قطعة جينومية، تظهر لنا صفحة البيانات كملف مستقل تترتب فيه المعلومات بنظام خاص. فمثلا عند البحث عن تتابعات جين المالتيز maltase gene فى أحد أنواع الخميرة Candida albicans فسيظهر لنا الملف المبين بشكل (٥-٢)، حيث تحدد بعض من

المعلومات الأساسية، ففي السطر الأول تحت عنوان locus يبين به الرقم الكودى للمدخل XM 714334 يليه حجم التتابع 1713 bp ثم نوع الجزيء mRNA في هذه الحالة وفي النهاية تاريخ الإدخال. وتحت definition تحدد بيانات اسم الجين والاسم العلمى للكائن المدروس، والتفاصيل التقسيمية لهذا الكائن يمكن ملاحظتها تحت عنوان organism. أما أهم المعلومات لقاعدة البيانات فتوجد تحت عنواني accession و version حيث يحدد الرقم الكودى للعينة وهو ثابت لا يتغير حتى لو تم تغير وتدفيق نفس التتابعات، لكن يظهر لنا رقم جديد للمدخل يسمى «المعرف» identifier وهو فى حالتنا gi 68473242 وهو لزيادة التحقق لكل قاعدة. وعادة توجد عناوين أخرى تحدد المصدر والمرجع الذى قدم المعلومات وأماكن نشرها ودرجة التحقق منها وغيرها من المعلومات، وتحت عنواني features و comment تتواجد معلومات اضافية عن الكائن والجين والحالات المشابهة، وطبيعة الجزيئات المدخلة للقاعدة وعادة تستعمل الإختصارات التالية لتمييزها:

con = constructed sequence,	est = expressed sequence tag,
gss = genome survey sequences,	htg = high throughput genomic,
new = new since last release,	patent = patent, and
sts = sequence tagged site.	

وفي النهاية ترتب تتابعات الأحماض الأمينية المترجمة من هذا التتابع وبعدها تتابعات النيوكليوتيدات لهذا الجين والتي قدرها فى هذه الحالة ١٧١٢ زوج من القواعد.

LOCUS	XM_714334	1713 bp	mRNA	linear	PLN 29-SEP-2005
DEFINITION	Candida albicans SC5314 maltase (CaO19_7668), mRNA.				
ACCESSION	XM_714334	XM_435011			
VERSION	XM_714334.1	GI:68473242			
KEYWORDS					
SOURCE	Candida albicans SC5314				
ORGANISM	<u>Candida albicans SC5314</u> Eukaryota; Fungi; Ascomycota; Saccharomycotina; Saccharomycetes; Saccharomycetales; mitosporic Saccharomycetales; Candida.				
REFERENCE	1 (bases 1 to 1713)				
AUTHORS	Jones, T., et al.				
TITLE	The diploid genome sequence of Candida albicans				
JOURNAL	Proc. Natl. Acad. Sci. U.S.A. 101 (19), 7329-7334 (2004)				
PUBMED	<u>15123810</u>				
REFERENCE	2 (bases 1 to 1713)				
AUTHORS	Jones, T. et al.				
TITLE	Direct Submission				
JOURNAL	Submitted (16-APR-2004) Stanford Genome Technology Center, 855 California Avenue, Palo Alto, CA 94304, USA				
COMMENT	PROVISIONAL <u>REFSEQ</u> : This record has not yet been subject to final NCBI review. This record is derived from an annotated genomic sequence (NW_139452).				
COMPLETENESS:	incomplete on both ends.				
FEATURES	Location/Qualifiers				
source	1..1713 /organism="Candida albicans SC5314" /mol_type="mRNA" /strain="SC5314" /db_xref="taxon:237561" /chromosome="R"				
gene	1..1713 /locus_tag="CaO19_7668"				

CDS	/db_xref="GeneID:3638946"
	1..1713
	/locus_tag="CaO19_7668"
	/note="gene whose transcription is induced by maltose and sucrose and repressed by glucose; one of two genes similar to P.angusta MAL1 and to seven maltase and maltase-like genes in S. cerevisiae"
	/codon_start=1
	/transl_table=12
	/product="maltase"
	/protein_id="XP_719427.1"
	/db_xref="GI:68473243"
	/db_xref="GeneID:3638946"
/translation=	
"MSEHKWWKEAVVYQIWPASYKDSNGDGV.....GNY KLVLTNVDKDS KDALSPYEARMYVVD"	
ORIGIN	
atgagtgaaac ccatacgagg ctagaatgta ttagttgattaa //	

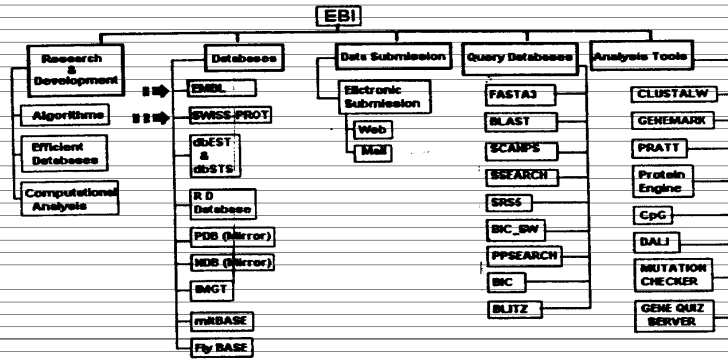
شكل (٢-٥) : صفحة من ملف في الـ GenBank توضح تنبؤات جين المالتيز في أحد أنواع الخميرة.

٢.٢.٥. معهد المعلوماتية الأوروبي EBI.

بدأ التفكير في إنشاء هذا المعهد سنة ١٩٨٦ وهو نتيجة التعاون بين العديد من الدول الأوروبية بغرض التخزين والتبويب الإلكتروني للمعلومات الوراثية خصوصا تنبؤات النيوكليوتيدات والبروتينات، وهو الإمتداد الطبيعي لعامل البيولوجيا الجزيئية الأوروبية European Molecular Biology Laboratory والمعروفة اختصارا EMBL والمنتشرة في العديد من الدول الأوروبية والمقر الرئيسي للمعهد يوجد في معمل البيولوجيا الجزيئية في بلدة Hinxton بإنجلترا. وأهداف المعهد مماثلة لأهداف المركز الأمريكي المناظر NCBI، ويمكن تلخيصها في الآتي:

- تقديم وتطوير تكنولوجيا تتبع المعلومات.
- تطوير وتحديث برامج دراسة تكنولوجيا المعلومات.
- توفير برامج تعليمية وتدريبية في مجالات تكنولوجيا المعلومات.

ولتحقيق تلك الأهداف قام المعهد بإنشاء العديد من الخدمات والمواقع على شبكة الأنترنت يمكن تلخيصها في شكل (٢-٥).



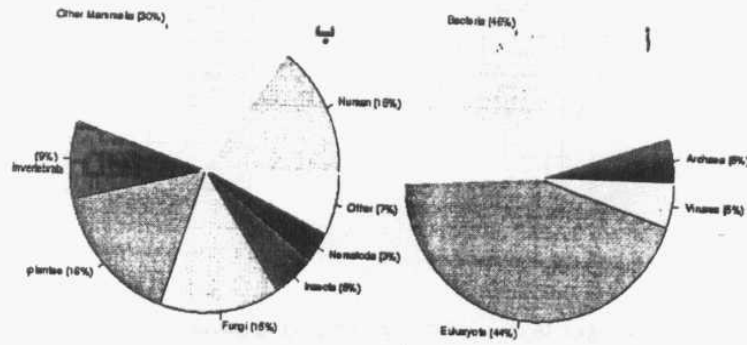
شكل (٢-٥) : رسم توضيحي يمثل المهام والخدمات التي يوفرها معهد المعلوماتية الأوربي EBI.

ومن جملة تلك المواقع والخدمات يبقى موقع EMBL الخاص بتتابعات النيوكليوتيدات والموقع SWISS-PROT لتتابعات البروتينات هما الأهم بالنسبة للوراثيين، والقاعدة EMBL لا تختلف كثيرا عن نظيرتها في المركز الأمريكي GenBank، لذا سنكتفي هنا بالتعريف بالقاعدة SWISS-PROT

١.٢.٢.٥. قاعدة معلومات البروتينات SWISS-PROT.

تشمل قاعدة البيانات المعروفة بأسم SWISS-PROT على البيانات الخاصة بتتابعات البروتينات وقد بدأت القاعدة من قسم الكيمياء الحيوية بجامعة جنيف Geneva University (المعروف الآن بأسم المعهد السويسري للمعلوماتية Swiss Institute of Bioinformatics) وقد انفرد الـ EBI بأصدار هذا الموقع مع مطلع التسعينات من القرن الماضي، وعنوانه على الأنترنت

هو <http://www.expasy.org/sprot/userman.html>. وفي عام ١٩٩٦ أضيفت خدمه جديدة لهذا الموقع عرفت بأسم TR-EMBL وهو تحت موقع يختص بتخزين بيانات ترجمة سلاسل النيوكليوتيدات المخزنة في EMBL أوتوماتيكيا بإستعمال برامج الكمبيوتر. وحتى سبتمبر من ٢٠٠٥ بلغ عدد المحفوظات accessions إلى 2151724 تتابع بروتيني موزعة على كافة الممالك الحية كما هو موضح بشكل (٥-٤).



شكل (٥-٤) : رسم بياني يمثل نسب تتابعات البروتينات المخزنة بموقع

SWISS-PROT / TR-EMBL : (أ) الممالك المختلفة (ب) الكائنات حقيقية النواة.

وقواعد التبويب وترتيب المعلومات في هذه القاعدة متماثل لحد كبير للقواعد والنظم المتبعة في القاعدة الأمريكية السابق بيانها، وشكل (٥-٥) يمثل صفحة لأحد المحفوظات بقاعدة SWISS-PROT، وهي بيانات حول البروتين الإنساني tyrosine protein-kinase والذي يتضح من تتابعاته أنه يتكون من ٥٣٦ حمضا أمينيا.

٣.٢.٥. البنك الياباني لبيانات الـ DNA (DDBJ).

البنك الياباني للمعلومات الخاصة بتتابعات الـ DNA هو ثالث القواعد الكبيرة والأكثر استعمالاً، وهو الآخر ترجع بداياته إلى سنة ١٩٨٦ حيث أنشئ كوحدة تابعة للمعهد القومي للوراثة NIG بمشيما بالقرب من طوكيو تحت إشراف وزارة التعليم والعلوم اليابانية. وما لبث أن انضم إلى التعاون الدولي في هذا المجال مع EBI الأوربي و NCBI الأمريكي. وفي سنة ٢٠٠١ تم الاعتراف بهذه القاعدة على أنها نواة مركز للمعلومات البيولوجية وبنك معلومات الـ DNA الذي يعرف اختصاراً بأسم (CIB-DDBJ). ويعتمد المركز على المدخلات التي يقدمها العلماء اليابانيون في هذا المجال وإن كان يقبل المعلومات من خارج اليابان أيضاً. والمعهد يدعم بشكل مؤكد البحوث في مجال البيولوجية الجزيئية والتتابعات الجينومية، وإن كان التعامل معه أكثر صعوبة لإعتماده على اللغة اليابانية.

شكل (5-5): عينة من أحد ملفات القاعدة SWISS-PROT.

٦- رصّ التتابعات

.Sequence Alignment

إعداد: أمير بسن

١.٦ مقدمة.

أيا ما كان تعريف الحياة فما هي إلا تعبير عن لغة من الأحماض النووية والبروتينات وغيرها من الجزيئات البيولوجية. وهي لغة تبلغ من الثبات حدا تبدو معه التباينات الهائلة بين شتى الكائنات الحية، الدقيقة منها والراقية والندشرة منها والباقية، وكأنها لا تعدو أكثر من مجرد لهجات محلية للغة أم واحدة. صحيح أننا نعلم اليوم أن ما يصدق لبكتيريا القولون إ. كولاى لا يصدق للفيث كما كان يظن عند بدء فك طلاسم هذه اللغة، إلا أننا نعلم كذلك أن هذين الكائنين على تباينهما حجما وتعقيدا وتطورا يتشاركان في جوهر واحد هو لغة الحياة الصماء في خلاياهما. فمن ألفباء من أربعة أحرف فقط تمثل وحدات البناء الكيميائية للبناء ينبثق بناء كامل من العمليات الحيوية يعتبر الإنسان أكثر تعبيراته تعقيدا. ويعد الكثُلف واستخدام تلك "الألفباء" في تشكيل "كلمات وجمل" جديدة هو مصب اهتمام حقل البيولوجيا الجزيئية الحديثة وأعظم تحدياته.

والحياة موجودة من قبلنا بأربعة بليون سنة وستظل موجودة بعدنا (إن لم نقض عليها نحن بسياساتنا المدمرة)، وهي ليست بحاجة لن يفك طلاسمها. ولكنها طوال فترة تواجدنا على الأرض وهي تشق طريقها بنفسها وتعيد صياغة الكوكب لتوائم ظروفه وجودها عليه، سائرة في ذلك في خط مستقيم نحو زيادة حجم وقدرات ذلك العضو الحسى العصبى المسمى بالمخ والذي بلغ أكثر تعقيد له في المخ البشرى، جاهلة كل شئ عن وجودها وجوهرها، حتى تسنى لهذا الأخير أن يحل اللفز ويقرأ كتابها.

وإنه ليحق لنا أفراد هذا النوع البشرى وأبناء هذا الجيل أن نفخر بأننا كنا أول من نجح في قراءة سفر الحياة في خلايا الكائنات الحية وفي خلايانا ذاتها، ويكفى أن القرن الماضى قد انتهى بإعلان مسودة الجينوم البشرى بعد خمسين عاما فقط من اكتشاف المادة الوراثية ومعرفة تركيبها، وهو ما سوف ينعكس بدوره على طبيعة وتوجهات البحوث البيولوجية في القرن الحادى والعشرين. فمنذ التسعينيات والجمع

العلمى يعيش فى ثورة تشبه تلك التى واكبت اكتشاف الذرة فى مطلع القرن العشرين، وكأى ثورة علمية تكمن عظمتها فى مقدار ما تطرح من أسئلة جديدة أكثر من مقدار ما تمنح من إجابات. لقد أدت علوم الذرة وميكانيكا الكم إلى قلب الفيزياء الكلاسيكية رأساً على عقب، فما الذى سوف يؤدى إليه قراءة سفر الحياة ؟ وإذا كانت الحيطه والحذر يتوخيان بنا ألا ننحرف خلف نبوءات واسعة قد يثبت كذبها فيما بعد، إلا أننا نتنبأ بل ونؤمن لشد الإيمان أن أكثر الاكتشافات إثارة فى مجالات البيولوجيا والتى سوف يكشف عنها الغد هى تلك التى لا تزال تندرج اليوم فى خانة الأمور غير المتوقعة أو حتى المستحيلة.

لقد كان من حسن حظ البيولوجيا الجزيئية أن واكبت ثورة أخرى فى مجال تكنولوجيا المعلومات، فإن الحجم المذهل للبيانات الجزيئية ولعماطها الخفية والماكرة أدت إلى عدم وجود مناص من استخدام قواعد بيانات وأدوات تحليل تعمل بواسطة الكمبيوتر. فالتحدى الأكبر إذن هو إيجاد مداخل جديدة للتعامل مع حجم البيانات وتعقيدها ولد الباحثين بطرق أسهل للوصول إلى المعلومة وأدوات الكمبيوتر حتى يتسنى لنا فهم أفضل ليراثنا الجينى ودوره فى كل من الصحة والمرض. وهو ما تصدى له علم المعلوماتية الحيوية. فالمعلوماتية الحيوية bioinformatics هى طريقة جديدة لقراءة البيولوجيا أكثر منها فرع من العلوم يستخدم الكمبيوتر فى الأبحاث البيولوجية كما يحلو للكثيرين النظر إليها (وهو ما يميزها عن الحوسبة الحيوية - biocomputing)، وهى تنظر للغة الحياة (تتابعات الدنا أو تتابعات البروتينات وأدبيات البيولوجيا) على أنها مجموعة من البيانات الخام data التى لا بد من معالجتها حتى يتم صياغة معنى لها، هذا المعنى هو المعلومة information التى هى غاية الباحث أى باحث. وتكون هذه المعالجة processing إما بتخزينها فى قواعد بيانات يسهل الوصول إليها، أو — وهذا هو الأهم — إيجاد الصلة بين هذه البيانات المختلفة أى ما هو التشابه بينها، وذلك فى ضوء محاور ثلاثة :

- يحدد تتابع الدنا تتابع البروتين
- يحدد تتابع البروتين تركيب البروتين (شكله الفراغى)
- يحدد تركيب البروتين وظيفة البروتين

وتجدر الإشارة إلى أن التشابه similarity في حد ذاته هو مشكلة انطولوجية، فحتى اليوم لا يمكن إيجاد مفهوم واحد للتشابه، فمثلا قد يتشابه نوعان من الأعمدة مثل الأعمدة الكورنثية والأعمدة الإيونية في العمارة الإغريقية في أن كلا منهما يتكون من بناء أسطوانى طويل ينتهى بإكليل مزخرف وهو تشابه تركيبى structural similarity يعكس انتماء كل منهما لنفس الفن، كما أنه قد يتشابه شيئان لا صلة تركيبية بينهما مثل معجون الأسنان وماكينه الحلاقة حيث يتشابهان في أن كلا منهما يلعب دورا في سيناريو الصباح اليومى وهو تشابه وظيفى functional similarity، كذلك قد يتشابه مقعد خشبى وثوب من الكتان في أن كلا منهما قد صنع من مادة خام نباتية وهو ما يعرف بالتشابه التطورى evolutionary similarity (أى التشابه من حيث آليات التخليق). وفي البيولوجيا يتشابه الصرصور مع الذبابة تشابها تركيبيا (يتكون الجسم من رأس وصدر وبطن ويحمل الصدر ستة أزواج من الأرجل وزوجين من الأجنحة بينما تحمل الرأس العيون وقرنى الاستشعار وأجزاء الفم)، بينما تتشابه أجنحة الطيور مع أجنحة الخفافيش تشابها وظيفيا (وإن اختلفتا تركيبا)، بينما يتشابه الكنجارو مع الفأر مثلا تشابها تطوريا (فكل منهما حيوان ثديى وإن اختلفت آليات الولادة في كل منهما تركيبا ووظيفة).

والخطوة الأولى لإيجاد أى تشابه بين البيانات الجزيئية تستخرج منه معلومة هي ترتيب التتابعات قيد الدراسة (سواء كانت تتابعات دنا أو تتابعات بروتين أو تتابع دنا مع تتابع بروتين) بصورة يسهل معها تقدير التشابه ومعرفة طبيعته (تركيبى أو وظيفى أو تطورى) وهو ما يعرف باسم رص التتابعات sequence alignment. وإذا كان كل من التركيب والوظيفة عادة ما يعكسان التطور في البيولوجيا، حتى أنه ليصعب فصله عنهما أو تفسيرهما بدونه، فإننا سوف نعمد بعد شرحنا لبعض أساسيات لغة الحياة إلى شرح التغيرات التطورية في التتابعات وطرق تقدير معدلات وانماط تلك التغيرات وذلك في معرض حديثنا عن حسابيات algorithms وبرمجيات software رص التتابعات والذي هو موضوع بحثنا هذا.

٢.٦. قواعد لغة الحياة.

١.٢.٦. تتابعات النيوكليوتيدات.

بخلاف بعض الفيروسات، فإن المادة الوراثية في كل الكائنات الحية يحملها الحامض النووي الديوكسي ريبوزي الذي يعرف اختصاراً بالدنا DNA، والذي يتكون من شريطين يلتفان حول بعضهما البعض مشكلين لولياً مزدوجاً يميني الالتفاف. ويتكون كل شريط منهما من سلسلة من أربع نيوكليوتيدات nucleotides اثنان منهما من نوع البيورين ويرمز لهما بالحرف R (وهما الأدينين A والجوانين G) واثنان من البيريميدين ويرمز لهما بالحرف Y (وهما الثايمين T والسيتوزين C). ويرتبط شريطا الدنا من خلال روابط هيدروجينية بين نيوكليوتيداهما حيث دائماً ما ترتبط نيوكليوتيدة من البيورين بنيوكليوتيدة خاصة من البيريميدين، وتسمى الرابطة بين الأدينين والثايمين باسم الرابطة الضعيفة ويرمز لها بالحرف W، بينما تسمى الرابطة بين الجوانين والسيتوزين باسم الرابطة القوية ويرمز لها بالحرف S، وتكتب الروابط بوضع نقطتين بين رموز النيوكليوتيدات مثل C:G أو A:T فيما يعرف باسم أزواج القواعد القانونية canonical base pairs. ويوضح الجدول (١-٦) أبجدية كتابة تتابعات النيوكليوتيدات.

ويحدد اتجاه الروابط الفوسفوديسترية التي تربط بين نيوكليوتيدات الشريط الواحد في السلسلة طبيعة الجزئ، وعليه فإن النيوكليوتيدات تكتب بترتيب النسخ أي في الاتجاه 3' → 5' فقط، وبالنسبة لموقع معين على التتابع فإن كلا من الطرفين 3' أو 5' يسميان الجانب العلوي upstream أو الجانب السفلي downstream على التوالي. ويسمى الشريط الذي يحتوي على أكثر من 50% من البيورينات باسم الشريط الثقيل heavy strand، بينما يسمى الشريط الذي يحتوي على أكثر من 50% من البيريميدينات باسم الشريط الخفيف light strand.

جدول (٦-١) - أبجدية الدنا.

الحرف	التفسير
A	أدينين
C	سيتوزين
T	ثايمين
G	جوانين
W	روابط ضعيفة (A,T)
S	روابط قوية (C,G)
R	بيورينات (A,G)
Y	بيريميدينات (C,T)
K	كينو (T,G)
M	أمينو (A,C)
B	C أو G أو T
D	A أو G أو T
H	A أو C أو T
V	A أو C أو G
N	A أو C أو G أو T
-	لا نيوكليوتيدة (فجوة)

أما الرنا RNA فإنه يتواجد إما في صورة شريط مفرد أو مزدوج، وهو يختلف عن الدنا في احتواءه على سكر ريبوز بدلا من سكر الديوكسي ريبوز وفي أن نيوكليوتيدة اليوراسيل U تحل محل الثايمين، كما أن أزواج القواعد القانونية فيه تشمل G:U كذلك (بينما لا توجد رابطة G:T في جزئ الدنا). وتسمى هذه النيوكليوتيدات مجتمعة (A، C، G، T، U) باسم النيوكليوتيدات القياسية standard nucleotides، ويلاحظ أن بعض جزيئات الرنا وخاصة الرنا الناقل t-RNAs قد تحتوي على بعض النيوكليوتيدات غير القياسية والتي تشتق من النيوكليوتيدات القياسية. ويقاس طول الحمض النووي المفرد (سواء دنا أو رنا) بعدد نيوكليوتيداته، بينما يقاس الجزء المزدوج بعدد أزواج القواعد (زق) (bp) base pairs.

وتسمى المادة الوراثية كلها في الكائن الحي باسم الجينوم genome، وهو ينقسم إلى جزء يحمل جينات ويسمى الدنا الجيني genic DNA وجزء لا يحمل جينات يسمى الدنا غير الجيني nongenic DNA (وقد يكون هذا الجين حاملا لجينات لم يتم الكشف عنها بعد). ويسمى كل ما يتم نسخه من الجينوم من رنات باسم الترانسكريبتوم transcriptome، بينما يسمى كل ما يتم تشفيره (ترجمته) من بروتينات باسم البروتيوم proteome.

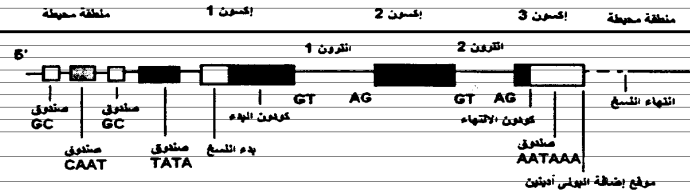
وبالرغم من أن تعريف الجين gene (وهو محور أبحاث علم الوراثة genetics عموما) لا يزال غامضا، إلا أن أكثر التعريفات الحديثة شيوعا هي تلك التي تعتبره قطعة من الدنا أو الرنا مسؤولة عن وظيفة معينة، وليس بالضرورة أن يقتضى تحقيق هذه الوظيفة ترجمة أو حتى نسخا. وثمة ثلاثة أنواع من الجينات :

- جينات مشفرة للبروتينات protein-coding genes، وتلك يتم نسخها إلى رنات ومن ثم ترجمتها إلى بروتينات.
- جينات مخصصة للرنا RNA-specifying genes، وتلك يتم نسخها فقط.
- جينات لا يتم نسخها untranscribed genes.

ويسمى النوعان الأولان باسم الجينات التركيبية structural genes (وأحيانا تطلق هذه التسمية على النوع الأول فقط منها).

وتتكون الجينات المشفرة للبروتينات في حقيقيات النواة (شكل ١-٦) من أجزاء يتم نسخها وأخرى لا يتم نسخها، والأجزاء التي لا يتم نسخها منطقتان محيطتان flanking regions، واحدة تقع عند الطرف 5' وتحتوى على تتابعات خاصة تسمى الإشارات signals تكون مسؤولة عن حث عملية النسخ وضبط إيقاعها وتوقيتها وتعيين النسيج الذى سوف يتم التعبير عنها فيه، وأحيانا تسمى هذه التتابعات باسم المحفزات promoters، وتسمى المنطقة التي تحتويها باسم منطقة الحفز، وهي تحتوى على الإشارات التالية مثل صندوق TATA الذى يقع على بعد ٢٧-١٩ زها على الجانب العلوى من نقطة بدء النسخ ويتحكم في اختيار هذه النقطة، وصندوق CAAT الذى يقع أكثر بعدا على الجانب العلوى، وصناديق CG المتباعدة عددا وموقعا حول صندوق CAAT والتي تشترك

معه في الارتباط ببوليميراز الرنا عند النسخ. وتجدر الإشارة إلى أن هذه التتابعات ليست بالضرورة موجودة في كل الجينات. أما الجانب السفلي فيحتوي على إشارات إنهاء عملية النسخ. ونتيجة لقصر معلوماتنا عن كل من إشارات المنطقتين فإننا لا نزال إلى الوقت الحاضر عاجزين عن تحديد نقطتي بدء وانتهاء الجين بدقة.



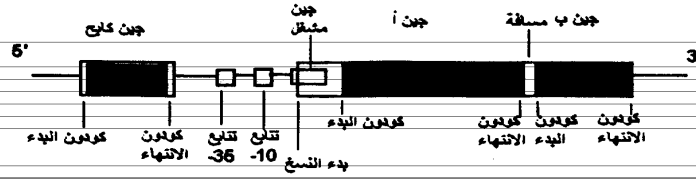
شكل (١-٦) – رسم تخطيطي لجين مشفر للبروتين نموذجي في حقيقيات النواة ومكوناته التركيبية ويشر الخط المتقطع إلى احتمالية ألا يتماثل موقع إضافة البولي أدينين مع موقع انتهاء النسخ.

ويسمى الرنا المنسوخ داخل النواة باسم الرنا النووي غير المتجانس hnRNA، كما يسمى شريط الدنا الذي ينسخ منه الرنا باسم شريط الاتجاه المضاد antisense strand بخلاف الشريط المقابل الذي لا يتم نسخه ويكون متماثلاً مع شريط الرنا والذي يسمى شريط الاتجاه الأصلي sense strand. ويبدأ النسخ عادة من الدنا عند نقطة تسمى موقع حث النسخ transcription initiation site (والتي تقابل قنسوة cap الرنا)، بينما ينتهي عند نقطة تسمى موقع انتهاء النسخ transcription termination site (والتي ليس بالضرورة أن تقابل موقع إضافة ذيل البولي أدينين على شريط الرنا الرسول mRNA الناضج فقد تقع بعد هذا الموقع). ويتكون بادئ الرنا الرسول pre-mRNA من إكسونات وانترونات، أما الانترونات introns فهي تلك التتابعات التي يتم قطعها والتخلص منها خلال عملية معالجة الرنا RNA processing، أما سائر التتابعات الأخرى والتي تبقى على طول شريط الرنا بعد عملية المعالجة ويتم الوصل بينها فتسمى إكسونات exons. وبالإضافة لعملية الوصل splicing يبلغ بادئ الرنا الرسول مرحلة

النضج بعد وضع قلنسوة من الجوانين الميثيل عند الطرف 5' وتحلل النيوكليوتيدات التي قد تقع بعد نقطة إضافة ذيل عديد الأدينين ثم إضافة الأخير (والذي قد يبلغ طوله ٢٠٠-١٠٠ قاعدة أدينين). كذلك تتعرض جزيئات الرنا للعديد من عمليات التحكم النسخي وما بعد النسخي transcriptional and post-transcriptional processes.

وتصنف الانترونات تبعا لطريقة الوصل إلى انترونات ذاتية الوصل self-splicing introns مثل انترونات جينات الميتوكوندريا والكلوروبلاست، وانترونات سبليسيوسومية spliceosomal introns مثل انترونات الجينات النووية والتي يتم انزعاجها من الرنا بواسطة معقد إنزيمي يدعى السبليسيوسوم spliceosome. وتسمى الأطراف 5' من الانترون بالمواقع المانحة donor sites وغالبا ما تكون GT، بينما تسمى الأطراف 3' بالمواقع المستقبلة acceptor sites وغالبا ما تكون AG فيما يعرف باسم قاعدة GT-AG. ويختلف عدد وتوزيع الانترونات من جين إلى آخر في حقيقيات النواة، وإن كانت معظم جينات الققازيات تتكون أساسا من انترونات

أما في أوليات النواة، فإن جيناتها (شكل ٢-٦) لا تحتوي على انترونات، ولها محفزان للنسخ على الجانب العلوي يقعان على مسافة ١٠ قواعد (صندوق Pribnow ويتكون من التتابع TATAAT) و٢٥ قاعدة (تتابع TTGACA أو أحد تنويعاته) مع احتمال وجود محفزات أخرى. كذلك تنتظم الجينات التركيبية في أوليات النواة على طول الجينوم مشكلة وحدة تعبير وراثي واحدة تنتج رنا رسولا طويلا يترجم إلى عدة بروتينات مختلفة، وتسمى هذه الجينات مجتمعة مع الجينات المنظمة لتعبيرها باسم الأوبرون operon.



شكل (٢-٦) — رسم تخطيطي لأوبيرون مستحث في أوليات النواة.

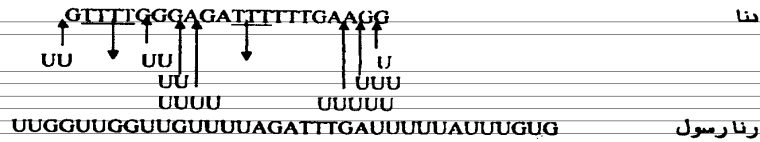
وسواء ترجم جزئ الرنا إلى بروتين أو لم يترجم فإنه عادة ما يعقب عملية

النسخ عدد من التحويلات تعرف باسم تحرير الرنا RNA editing. وقد تبلغ عملية

التحويل هذه حدا يصعب معه إيجاد التشابه بين جزئ الرنا وجزئ الدنا الذي نسخ

منه، عندئذ يسمى غالب الجين باسم كريبتوجين (جين خفي) cryptogene

(شكل ٢-٦).



شكل (٢-٦) — مقارنة بين منطقة ممثلة لشريط الاتجاه الصحيح لجين الميتوكوندريا المشفر للوحدة

الثانوية III لأوكسيداز السيتوكروم C في التريبانوسوما بروسي (Trypanosoma brucei) مع رنا

رسول بعد التحرير. تشير الأسهم إلى أعلى إلى إدراج U بينما تشير الأسهم إلى أسفل إلى حذف T

[من Feagin et al. (1988)].

كذلك يوجد عدد من الجينات لا تتم ترجمتها ولا حتى نسخها بيد أنها

يكون لها وظائف حيوية حيث أنها تعمل كمواقع ارتباط للإنزيمات في عمليات

تضاعف الدنا (جينات التضاعف replicator genes) والعبور الوراثي (جينات العبور

الوراثي recombinator genes) والتتابعات التيلوميرية telomeric sequences

والانقسام الخلوي (جينات الانعزال segregator genes) وتتابعات مواقع الارتباط

attachment sites بالبروتينات الهيكلية والإنزيمات والهرمونات والنواتج الأيضية وتتابعات المواقع البنائية constructional sites والتي لها علاقة بشكل الكروموسومات.

وقد تظهر أيضا بعض قطع الدنا غير الجينية تشابها شديدا لجينات عاملة ولكنها تحتوى على عيوب (طفرات) تمنع التعبير عنها بصورة صحيحة وتسمى فى هذه الحالة جينات كاذبة pseudogenes. وعادة ما تكتب بوضع الرمز ψ قبل اسم الجين العامل المشابه لها فى معظم الأدبيات وإن كانت فى قواعد بيانات الجينات الموجودة على الكمبيوتر يستبدل هذا الرمز بالحرف P.

٦.٢.٢. تتابعات البروتينات.

تتكون البروتينات proteins فى كل الكائنات الحية من ٢٠ حمضا أمينيا amino acids اوليا مذكورة ابجديتها فى الجدول (٦-٢). ويتكون كل حمض أمينى من مجموعة أمين NH_2 - ومجموعة كربوكسيل COOH - على جانبى ذرة كربون مركزية تسمى ألفا carbon α يرتبط بها كذلك ذرة هيدروجين ومجموعة السلسلة الجانبية R group -، وتكون الأخيرة هى المسئولة عن تمييز حمض أمينى عن الآخر، حيث أنها تتباين فى الحجم والشكل والشحنة والقدرة على تشكيل الروابط الهيدروجينية والتركيب والتفاعلية الكيميائية، ومن ثم فإن تصنيف البروتينات يتوقف على خواص سلاسلها الجانبية (شكل ٦-٤).

جدول (٢-٦) - أبجدية الأحماض الأمينية.

الحرف	الحمض	الحرف	الحمض
D	حمض الأسبرتك	R	أرجينين
E	حمض لجلوتامك	N	أسبرجين
C	سستين	A	الالين
S	سيرين	I	ايزوليوسين
V	فالين	P	برولين
F	فينيل الالين	W	تربتوفان
K	ليسين	Y	تيروسين
L	ليوسين	T	ثريونين
M	ميثيونين	Q	جلوتامين
H	هستيدين	G	جليسين

ويتكون البروتين من عدد من السلاسل البوليببتيدية لكل منها تتابعها الخاص

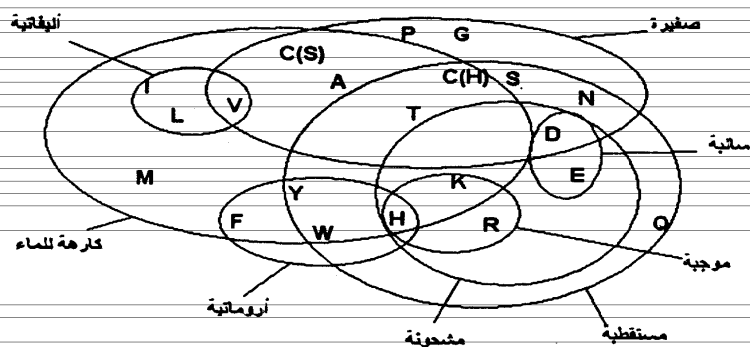
من الأحماض الأمينية التي ترتبط مع بعضها البعض بواسطة روابط ببتيدية peptide

bonds. ويسمى كل حمض أميني في السلسلة البوليببتيدية باسم فضالة residue، وتقرأ

السلسلة من الطرف الأميني أو النهاية النيتروجينية amino or N terminus إلى

الطرف الكربوكسيلي carboxyl or C terminus. وكما سبق أن اشرنا فإن تتابع الأحماض

الأمينية في السلاسل البوليببتيدية للبروتينات هو الذي يحدد تركيب هذه البروتينات.



شكل (٤-٦) - رسم فن يبرز تقسيم الأحماض الأمينية الأولية العشرين إلى مجموعات متداخلة تبعا للحجم وتركيب السلسلة الجانبية والاستقطاب والشحنة والكارهة للماء. لاحظ وقوع السيستين في مجموعتين كسستين C(H) وكسستين C(S).

وثمة أربعة مستويات لتركيب البروتين، فالتركيب الأولي primary structure

هو ببساطة الترتيب الخطي لفضالات الأحماض الأمينية على طول تتابع البوليببتيد. أما

التركيب الثانوي secondary structure فيعني الترتيب الفراغي أو التثنى folding

لفضالات الأحماض الأمينية المتجاورة في التركيب الأولي، ولعل أشهر أنواع هذا التركيب

اثنين هما اللولب ألفا α helix وهو تركيب عصوي تنتظم فيه الأحماض الأمينية في

شكل حلزون يميني الالتفاف تكون فيه سلاسلها الجانبية كلها للخارج وتنتظم الروابط

الهيدروجينية بين مجاميع الأمينو ومجاميع الكربوكسيل على مسافة كل أربعة أحماض

أمينية في التتابع، وتركيب صفيحة بيتا β sheet وفيه تنتظم السلاسل المتوازية من

الأحماض الأمينية بواسطة روابط هيدروجينية بين السلاسل المتجاورة. كذلك فبعض

مناطق البروتين لا ينتظمها أي تركيب ثانوي ثابت وتظل في التفاف عشوائي

random coil. وعلى العموم فإن التركيب الثانوي لمعظم البروتينات يكون مجموعة من

اللولب ألفا والصفائح بيتا والالتفافات العشوائية.

ويسمى التركيب الفراغى ثلاثى الأبعاد للبروتين باسم التركيب الثلاثى tertiary

structure ويعنى الترتيب الفراغى لفضالات الأحماض الأمينية غير المتجاورة فى التركيب الأولى، حيث ترتبط التراكيب الثانوية بواسطة قوى تساهمية وقوى غير تساهمية مثل الروابط الهيدروجينية والتفاعلات غير المحبة للماء وجسور الأملاح بين الفضالات مضادة الشحنة وروابط الديكربيتيد بين السستيينات. وفى هذا التركيب تعتمد الأحماض الأمينية ذات السلاسل الجانبية الكارهة للماء إلى الاندفاع فى قلب البروتين، بينما تنتظم الأحماض الأمينية ذات السلاسل الجانبية المحبة للماء على سطح البروتين.

وقد يتكون البروتين من أكثر من سلسلة بوليبيبتيدية واحدة تدعى كل منها فى

هذه الحالة بوحدة ثانوية subunit، ويقال عن البروتين إن له تركيبا رباعيا

quaternary structure وهو ما يعنى الترتيب الفراغى لهذه الوحدات الثانوية وطبيعة

ارتباطها ببعضها البعض. وإلى اليوم، يعتبر التنبؤ بتركيب البروتين من تتابعه فقط أحد أكبر تحديات الكيمياء الحيوية والمعلوماتية الحيوية نتيجة لقصر معلوماتنا عن طبيعة تكوين الروابط.

وتحتاج بعض البروتينات إلى بعض المركبات غير البروتينية للقيام بوظيفتها،

وتسمى تلك المواد باسم المجموع الترفيعية prosthetic groups، ويطلق على معقد

البروتين مع هذه المواد الترفيعية اسم البروتين الكامل holoprotein، بينما يسمى

البروتين بدون هذه المواد باسم البروتين الناقص apoprotein.

٢.٢.٦. الشفرة الوراثية.

اثناء عملية تخليق البروتين تتم ترجع تتابع الرنا الرسول إلى سلسلة

بوليببتيدية بواسطة الرنات الناقلة tRNAs، حيث يكون لكل من الأحماض الأمينية

العشرين رناه الناقل الخاص به والذي يحمل الكودون المضاد anticodon لثلاثى

النوكليوتيدات (الكودون codon) المحمول على الرنا الرسول. وتبدأ الترجمة عند كودون

البداة initiation codon وتنتهى عند كودون الانتهاء termination codon، وتسمى

المرحلة التى تصل إليها الترجمة على التتابع باسم إطار القراءة frame reading .

وما أن تنتهى الترجمة، حتى يتم تنشيط البروتين وظيفيا إثر بعض العمليات بعد النسخية مثل تحويل بعض الأحماض الأمينية الأولية أو نزع التتابعات الطرفية أو النقل داخل أو بين الخلايا أو إضافة المجاميع الترفيعية أو وصل البروتينات.

وكأى لغة فلفلة الحياة قواعد تنظم العلاقة بين تتابعات الدنا والبروتينات تسمى الشفرة الوراثية genetic code (جدول ٢-٦). وبخلاف بعض الاستثناءات، فإن الشفرة الوراثية كونية universal (أى تشترك فيها كل الكائنات حقيقية وأولية النواة)، غير أن بعض العلماء يفضلون تسميتها بالشفرة الوراثية القياسية نتيجة لوجود هذه الاستثناءات. ويلاحظ أن بعض الأحماض الأمينية يشفر لها أكثر من كودون واحد فى ظاهرة تعرف باسم التنكسية degeneracy وفيها تسمى هذه الكودونات التى تشفر لنفس الحمض الأميني باسم الكودونات المترادفة synonymous codons. كذلك فإن الكودونات المترادفة التى تختلف فقط فى النيوكليوتيدة الثالثة تسمى عائلة كودونية codon family.

وقد تغيب بعض الكودونات فى الجينات المشفرة للبروتينات تماما عن تتابع البروتين وتسمى كودونات غائبة absent codons، وأحيانا يعزى ذلك لغياب الرنا الناقل الخاص بها وتسمى الكودونات فى هذه الحالة باسم الكودونات غير المعينة unassigned codons ويمكن الكشف عنها بتوقف عملية الترجمة عند إطار قراءتها وعدم انفصال الرنا الرسول عن الريبوسوم.

٤.٢.٦. الطفرات.

تعرف الطفرات mutations بأنها أى تغير ينشأ فى تتابع الدنا نتيجة عادة للأخطاء التى تقع خلال التضاعف، وهى بذلك تعد المصدر الوحيد للاختلافات والمستحدثات فى التطور. حيث أن التطور evolution فى أبسط تعريفاته هو "الانحدار مع التحور" "descent with modification"، فإذا ما كانت هذه التحورات أو الاختلافات داخل نفس الأسرة (بين الآباء والأبناء) أو بين العشائر والأعراق داخل نفس النوع سمي

بالتطور الصغير microevolution، أما إذا تعدت هذه الاختلافات مستوى النوع وأدت إلى نشوء أنواع جديدة (وذلك بوضع حواجز تناسلية بين الأفراد الطافرة) فإنه يسمى بالتطور الكبير macroevolution. ومن الناحية التطورية تعتبر

جدول (٢-٦) – الشفرة الوراثية الكونية (وتظهر الاختلافات الموجودة في الشفرة الوراثية ليتوكوندرها الفقاريات بين هوسين).

الحمض الأميني	الكودون	الحمض الأميني	الكودون	الحمض الأميني	الكودون	الحمض الأميني	الكودون
C	UGU	W	UAU	S	UCU	F	UUU
C	UGC	W	UAC	S	UCC	F	UUC
Stop	UGA	Stop	UAA	S	UCA	L	UUA
(W)	UGG	Stop	UAG	S	UGG	L	UUG
W							
R	CGU	H	CAU	P	CCU	L	CUU
R	CGC	H	CAC	P	CCC	L	CUC
R	CGA	Q	CAA	P	CCA	L	CUA
R	CGG	Q	CAG	P	CCG	L	CUG
S	AGU	N	AAU	T	ACU	I	AUU
S	AGC	N	AAC	T	ACC	I	AUC
R (Stop)	AGA	K	AAA	T	ACA	I (M)	AUA
R (Stop)	AGG	K	AAG	T	ACG	M	AUG
G	GGU	N	GAU	A	GCU	V	GUU
G	GGC	N	GAC	A	GCC	V	GUC
G	GGA	E	GAA	A	GCA	V	GUA
G	GGG	E	GAG	A	GCG	V	GUG

الطفرات التي تحدث في الأنسجة التناسلية هي الأكثر أهمية حيث أنها تلك التي يتم توريثها للأجيال التالية، غير أنه من الناحية الطبية فإن للطفرات الجسمية أهمية كبرى كذلك إذ أنها قد تكون مسئولة عن بعض الأمراض مثل السرطان مثلاً.

وتنقسم الطفرات إلى :

١. طفرات إحلالية substitution mutations ، وفيها يتم إحلال

نيوكليوتيدة محل أخرى، فإذا حلت نيوكليوتيدة بيورينية محل نيوكليوتيدة أخرى من نفس النوع أو حلت نيوكليوتيدة بيريميدينية محل نيوكليوتيدة أخرى من نفس النوع سمي الإحلال transition (وهو أربعة أنواع : $A \rightarrow G$, $G \rightarrow A$, $C \rightarrow T$, $T \rightarrow C$)، أما إذا حلت نيوكليوتيدة بيورينية أو بيريميدينية محل نيوكليوتيدة أخرى من النوع الآخر سمي

الإحلال transversion (وهو ثمانية أنواع ، $A \rightarrow C, A \rightarrow T, C \rightarrow A, C \rightarrow G, G \rightarrow C, G \rightarrow T, T \rightarrow A, T \rightarrow G$). وقد تترجم هذه الطفرات على مستوى تتابع البروتين فإذا لم تحدث تغيراً سمين طفرات مترادفة synonymous، أما إذا أحدثت فإنها تسمى طفرات غير مترادفة nonsynonymous. وفي هذه الحالة إذا تم استبدال حمض أميني بآخر على تتابع البروتين فإنه يسمى استبدالاً replacement. ويجدر الإشارة إلى أنه ليست كل الطفرات المترادفة تكون صامتة silent أي لا تغير من البروتين المنتج، وإنما قد يؤثر بعضها على عملية الترجمة كتحويل إكسون إلى إنترون وما إلى ذلك. كذلك تصنف الطفرات غير المترادفة إلى طفرات خاطئة missense تغير حمضاً أمينياً بآخر، وأخرى لا معنى لها nonsense وهي تلك التي تحول كودونا مشفراً لحمض أميني إلى كودون انتهاء الترجمة ومن ثم ينتج البروتين ناقصاً.

٢. طفرات إعادة التوليفات الوراثية recombinations ، وتحدث نتيجة للعبور الوراثي crossing over أو نتيجة للتجور الجيني gene conversion.

٣. طفرات الحذف والإدراج (Indels) deletions and insertions ، وتحدث نتيجة العبور الوراثي غير المتساوي unequal crossing over والانتقال الوراثي DNA transposition والانتقال الأفقي للجينات horizontal gene transfer، وقد تؤدي هذه الطفرات إلى تغيير إطار القراءة ومن ثم تركيب البروتين كله، عندئذ تسمى طفرات نقل الإطار frameshift mutations.

٤. طفرات الانقلاب inversions، وتحدث نتيجة لكسور في الكروموسومات.

وتجدر الإشارة إلى أن الطفرات تختلف في معدلاتها بين الأنواع، فمثلاً يبلغ معدل الطفرات في فيروس الإنفلونزا ٢ مليون مرة أكثر من معدل الطفرات في الدنا النووي للفقاريات. كما تختلف الطفرات في توزيعها على طول الجينوم حيث تعرف المناطق ذات المعدلات العالية بالنقاط الساخنة hotspots مثل الجزر الغنية بالسيتوزين والجوانين CpG isochors. كذلك تختلف في أمتاطها فنجد أن transitions أكثر شيوعاً من

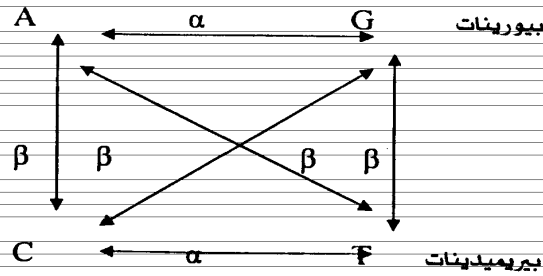
transversions، كذلك فإن النيوكليوتيدتين C وG أكثر تطفرا من النيوكليوتيدتين A وT.

ومن المعروف أن الطفرات تحدث "عشوائيا"، ولا يعنى هذا أن عشوائية الطفرات ترتبط بموقع الطفرة على الجينوم أو جميع أنواع الطفرات لها نفس التكرار، وإنما يرتبط مفهوم عشوائية الطفرة بتأثيرها على مواءمة الفرد حاملها. أى أن لى طفرة نفس احتمال الحدوث سواء كانت مفيدة أو ضارة، وهو ما عده ديزانسكى "نقصا فى الطبيعة" (Dobzhansky, 1970).

٢.٦. التغيرات التطورية فى التتابعات.

١.٢.٦. نموذج كيميورا ذو المقياسين للتغير فى تتابع النيوكليوتيدات.

لا يولد تتابع من عدم وإنما يعزى أى اختلاف بين تتابعين أو أكثر نتيجة للطفرات التى حدثت لكل منهما منذ زمن انحدارهما عن تتابعهما السلفى المشترك. ولدراسة ديناميات الإحلالات النيوكليوتيدية فلا بد من وضع عدة فرضيات بالنسبة لاحتمال إحلال نيوكليوتيدة محل أخرى. ولقد تم اقتراح عدة نماذج رياضية (را، ١٩٨١، ولكننا سوف نقصر هنا مناقشتنا على أحد أبسط هذه النماذج وهو النموذج ذى المقياسين لكيميورا (1980) Kimura's two-parameter model، الذى يقوم على افتراض أن معدل الإحلالات من نوع الـ transition ويرمز له بالرمز α لا يتساوى مع معدل الإحلالات من نوع الـ transversion ويرمز له بالرمز β (شكل ٥-٦).



شكل (٥-٦) - نموذج الإحلال النيوكليوتيدى ذى المقياسين، حيث قد لا يتساوى معدل الإحلالات من نوع

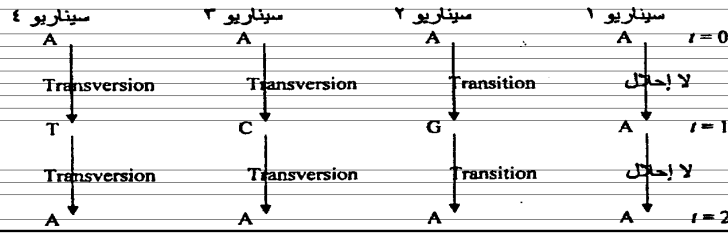
الـ transition ورمزه α مع معدل الإحلالات من نوع الـ transversion ورمزه β .

وعليه فإنه إذا ما كان لدينا موقع يحمل النيوكليوتيدة A عند الزمن $t = 0$ ، فإنه بعد وحدة زمنية واحدة يكون احتمال أن تتغير A إلى G هو α واحتمال أن تتغير إلى C أو T هو 2β ، وهكذا يكون احتمال أن يظل هذا الموقع يحمل النيوكليوتيدة A بعد وحدة زمنية واحدة هو

$$(1) \dots\dots\dots P_{AA(1)} = 1 - \alpha - 2\beta$$

أما عند $t = 2$ ، فهمة أربعة احتمالات : (١) ألا تتغير A في الوجدتين الزمنيتين، (٢) أن تتغير A إلى G في $t = 1$ ثم ترجع G إلى A في $t = 2$ ، (٣) أن تتغير A إلى C في $t = 1$ ثم ترجع C إلى A في $t = 2$ ، (٤) أن تتغير A إلى T في $t = 1$ ثم ترجع T إلى A في $t = 2$ ، (شكل ١-٦)، وعليه فإن

$$(2) \dots P_{AA(2)} = (1 - \alpha - 2\beta)P_{AA(1)} + \beta P_{TA(1)} + \beta P_{CA(1)} + \alpha P_{GA(1)}$$



شكل (١-٦) — أربعة سيناريوهات للحصول على A عند $t = 2$ تبعا لنموذج كيميوراى القياسين إذا كانت A موجودة على نفس الموقع عند $t = 0$.

وهو ما يمكن تعميمه إلى:

$$(3) \dots P_{AA(t+1)} = (1 - \alpha - 2\beta)P_{AA(t)} + \beta P_{TA(t)} + \beta P_{CA(t)} + \alpha P_{GA(t)}$$

وبالتفاضل في التغير الزمنى

$$(4) \dots \delta P_{AA(t)} / \delta t = -(\alpha + 2\beta)P_{AA(t)} + \beta P_{TA(t)} + \beta P_{CA(t)} + \alpha P_{GA(t)}$$

وبالمثل يمكننا الحصول على معادلات احتمالات الإحالات الثلاثة، والتي منها

نصل إلى

$$(*) \dots\dots\dots P_{AA}(t) = \frac{1}{4} + \frac{1}{4} e^{-4\beta t} + \frac{1}{2} e^{-2(\alpha+\beta)t}$$

أي أنه عند الاتزان ($t = \infty$)، يكون احتمال أي عدم إحلال $X(t)$ يساوي $\frac{1}{4}$.

وبعامة فإنه عند الزمن t ، يكون احتمال الإحلال من نوع الـ transition

$$(\dagger) \dots\dots\dots Y(t) = \frac{1}{4} + \frac{1}{4} e^{-4\beta t} - \frac{1}{2} e^{-2(\alpha+\beta)t}$$

بينما يكون احتمال الإحلال من نوع الـ transversion

$$(\ddagger) \dots\dots\dots Z(t) = \frac{1}{4} - \frac{1}{4} e^{-4\beta t}$$

حيث أن لكل نيوكليوتيدة نوع واحد من الإحلالات من نوع الـ transition ونوعين

من نوع الـ transversion K، وبالتالي فإن

$$X(t) + Y(t) + 2Z(t) = 1$$

٢.٤.٦. عدد الإحلالات النيوكليوتيدية بين تتابعين للدنا.

يعتبر عدد الإحلالات النيوكليوتيدية بين تتابعين هو أصل وأكثر المقاييس

شيوفا لتقدير مدى الاختلاف بينهما. فإذا اختلف تتابعان طول كل منهما N في عدد n

من المواقع، فإن النسبة n/N تعرف باسم درجة التشعب degree of divergence بينهما

أو مسافة هامنج Hamming distance، وعادة ما تكتب كنسبة مئوية. ويعيب هذه

الطريقة عدم قدرتها على الكشف عن الإحلالات المتعددة multiple substitutions or

multiple hits في نفس الموقع، كأن تتغير A إلى C ثم إلى T في تتابع وتتغير إلى T مباشرة

في التتابع الآخر فلا يكشف عنها بالرغم من وجود ثلاثة إحلالات بين هذين التتابعين

(شكل ٧-٦).

وهكذا يفضل التعبير عن مقدار هذا الاختلاف باستخدام عدد الإحلالات

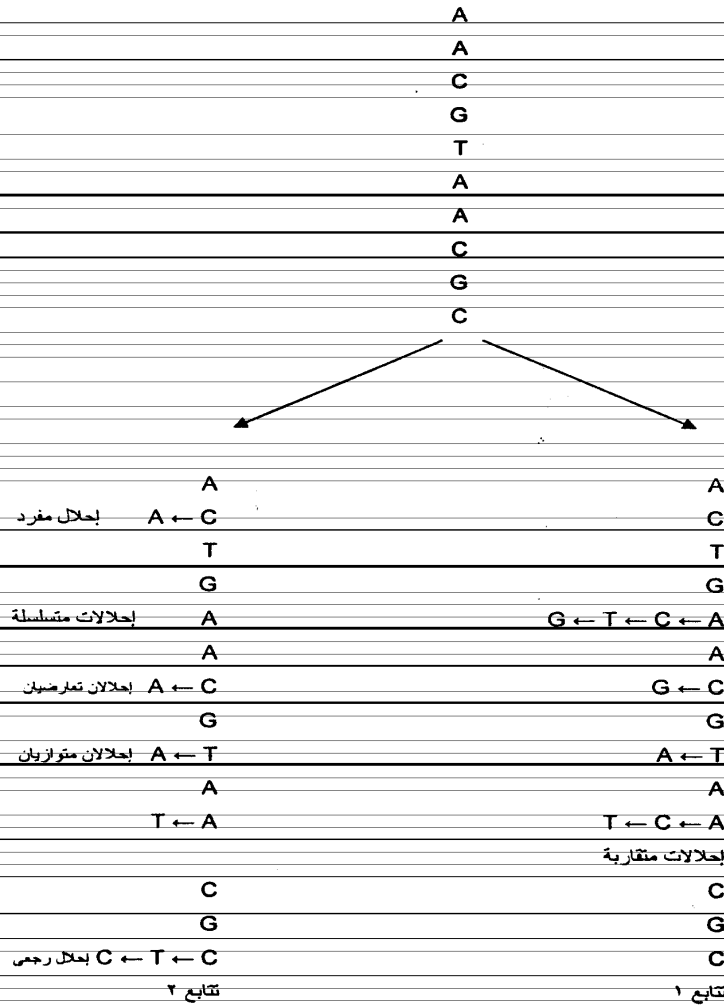
النيوكليوتيدية لكل موقع نيوكليوتيدى

number of nucleotide substitutions per nucleotide site ويرمز لها

بالرمز K بدلا من عدد الإحلالات النيوكليوتيدية الكلية بين التتابعين، خاصة إذا

اختلف التتابعان في أطوالهما. وتختلف طريقة حسابها بين إذا ما كان التتابع مشفرا

أو غير مشفر للبروتين نتيجة لاختلاف معدلات الإحلال بينهما.



شكل (٧-٦) - لتابعنا دنا منحدران من سلف مشترك وقد تراكمت الطفرات في كل منهما منذ زمن تشعبهما. لاحظ حدوث ١٢ طفرة بـعدة أنواع من الإحلالات حتى إن ٢ منها فقط هي التي يمكن مشاهدتها.

ففى حالة التتابعات غير المشفرة للبروتينات، وبافتراض أن عدد المواقع المقارنة بين التتابعين يساوى L ونسبة الاختلافات من نوعى $transition$ و $transversion$ يساويان P و Q على التوالى، فإن K تساوى

$$K = \frac{1}{2} \ln(1/1 - 2P - Q) + \frac{1}{2} \ln(1/1 - 2Q) \quad \dots \dots \dots (٨)$$

وذلك بتباين معاينة يساوى

$$V(K) = 1/L [P(1/1 - 2P - Q)^2 + Q(1/2 - 4P - 2Q + 1/2 - 4Q)^2 - ((P/1 - 2P - Q) + (Q/2 - 4P - 2Q) + (Q/2 - 4Q))^2] \quad \dots \dots \dots (٩)$$

ويلاحظ أننا قد استخدمنا هنا نموذج كيميورا ذا القياسين، وأن هذه المعادلات سوف تتغير إذا ما كنا قد استخدمنا نماذج ذات مقاييس أخرى.

أما بالنسبة للتتابعات المشفرة للبروتينات، فتعتبر عملية تحديد عدد الإحلالات النيوكليوتيدية أكثر تعقيداً من حالة التتابعات غير المشفرة للبروتينات حيث يجب التمييز بين الإحلالات المترادفة والإحلالات غير المترادفة.

وعند مقارنة تتابعين مشفرين للبروتين يتم إقصاء كودونى البدء والانتهاى حيث أن هذين الكودونين نادرا ما يتغيران مع الوقت.

أما عند حساب عدد الإحلالات لكل موقع لكل من الإحلالات المترادفة وغير المترادفة على حدة فلا بد أولاً من إيجاد القاسمين المناسبين، أى عدد المواقع المترادفة وغير المترادفة على التوالى. وهو ما يصعب عمله لسببين ؛
(١) يتغير تصنيف أى موقع مع الزمن، فمثلاً يعتبر الموقع الثالث للكودون CGG المشفر للأرجينين مترادفاً، ولكن إذا تغير الموقع الأول إلى T فإن الموقع الثالث والذي يصبح للكودون TGG المشفر للتربتوفان يصير غير مترادف.

(٢) بعض المواقع ليست مترادفة تماما أو غير مترادفة تماما، فمثلا إذا حدث إحلال من نوع transition في الموقع الثالث للكودون GAT المشفر لحمض الأسيرتيك فإنه سوف يكون مترادفا، أما إذا حدث إحلال من نوع الـ transversion فإنه سوف يكون غير مترادف، مع الأخذ في الاعتبار أن الإحلالات من نوع transition أكثر تكرارا من نوع transversion.

وفي سبيل ذلك تم وضع عدة نماذج سوف نتناول أحدها وهو النموذج الذي وضعه (Li et al. (1985 حيث أنه الأكثر شيوعا، وفيه يتم تقسيم المواقع النيوكليوتيدية إلى ثلاث فئات :

- موقع عديم التنكسية nondegenerate، إذا ما كانت جميع التغيرات عند هذا الموقع غير مترادفة، مثل الموقعين الأولين للكودون TTT المشفر للفينيل الانين.
- موقع ثنائي التنكسية twofold degenerate، إذا ما كان أحد التغيرات الثلاث عند هذا الموقع مترادفا، مثل الموقع الثالث للكودون TTT.
- موقع رباعي التنكسية fourfold degenerate، إذا ما كانت جميع التغيرات عند هذا الموقع مترادفة، مثل الموقع الثالث من الكودون GTT المشفر للفالين.

وعليه فإنه عند مقارنة تتابعين واتباع القواعد السابقة فإنه يتم عد أنواع المواقع الثلاثة لكل تتابع ثم حساب متوسط كل منها حيث يرمز للمواقع عديمة التنكسية بالرمز L_0 ، وللمواقع ثنائية التنكسية L_2 ، وللمواقع رباعية التنكسية L_4 . ثم يتم حساب عدد الإحلالات لكل نوع من المواقع على حدة، وذلك بعد تصنيف الاختلافات النيوكليوتيدية داخل كل فئة إلى اختلافات من نوع الـ transition ويرمز لها S_i واختلافات من نوع الـ transversion ويرمز لها بالرمز V_i حيث $i = 0, 2 \text{ or } 4$ تبعاً لدرجة التنكسية. ويلاحظ أنه بالتعريف فإن أي إحلال في موقع عديم التنكسية يكون غير مترادف، بينما أي إحلال في موقع رباعي التنكسية يكون مترادفا، ومن هنا تكمن المشكلة في المواقع ثنائية التنكسية حيث تكون جميع الإحلالات من نوع الـ transition مترادفة، بينما تكون جميع الإحلالات الأخرى من نوع الـ transversion غير مترادفة. كذلك لا بد من الانتباه إلى الشفرة الوراثية المستخدمة، حيث قد تختلف تنكسية الكودون من الشفرة الكونية إلى شفرة ميتوكوندريا الفقاريات مثلا، حيث أنه في الأخيرة لا توجد

استثناءات، بينما يوجد استثناءان في الشفرة الكونية وهما الموقع الأول للكودونات الأربعة المشفرة للأرجينين (CGA, CGG, AGA, AGG) والموقع الثالث للكودونات الثلاثة المشفرة للأيزوليوسين (ATT, ATC, ATA).
وتحسب نسبة الاختلافات من نوع transition عند مواقع ذات تنكسية من الدرجة i بين تتابعين من المعادلة

$$(10) \dots\dots\dots P_i = S_i / L_i$$

وبالمثل، تحسب نسبة الاختلافات من نوع transversion عند مواقع ذات تنكسية من الدرجة i بين تتابعين من المعادلة

$$(11) \dots\dots\dots Q_i = V_i / L_i$$

وباستخدام نموذج كيميورا ذي المقياسين لتقدير عدد الإحالات من نوعي transition ويرمز لها بالرمز A_i والـ transversion ويرمز لها بالرمز B_i من المعادلتين

$$(12) \dots\dots\dots A_i = \frac{1}{2} \ln(a_i) - \frac{1}{4} \ln(b_i)$$

و

$$(13) \dots\dots\dots B_i = \frac{1}{2} \ln(b_i)$$

وذلك بتباين

$$(14) \dots\dots\dots V(A_i) = [a_i^2 P_i + c_i^2 Q_i - (a_i P_i + c_i Q_i)^2] / L_i$$

و

$$(15) \dots\dots\dots V(B_i) = b_i^2 Q_i (1 - Q_i) / L_i$$

حيث L_i هو عدد المواقع ذات التنكسية من الدرجة i ، و $a_i = 1 / (1 - 2P_i - Q_i)$ ، و $b_i = 1 / (1 - 2Q_i)$ ، ويكون العدد الكلي للإحالات لكل موقع ذي تنكسية من الدرجة i ويرمز له بالرمز K_i وهو

$$(16) \dots\dots\dots K_i = A_i + B_i$$

بتباين

$$V(K_i) = [a_i^2 P_i + d_i^2 Q_i - (a_i P_i + d_i Q_i)^2] / L_i \quad (١٧)$$

حيث $d_i = b_i + c_i$.

ومما سبق وبعد التصحيح الذى وضعه Li و Pamilo and Bianchi (1993)

(1993) يكون عدد الإحالات المترادفة لكل موقع مترادف

number of synonymous substitutions per synonymous site ورمزه K_s

يساوى

$$K_s = [L_2 A_2 + L_4 A_4 / L_2 + L_4] + B_4 \quad (١٨)$$

بتباين

$$V(K_s) = [L_2^2 V(A_2) + L_4^2 V(A_4) / L_2 + L_4] + V(B_4) - [2b_4 Q_4 (a_4 P_4 - c_4 (1 - Q_4))] / L_2$$

$$+ L_4 \quad (١٩)$$

بينما يكون عدد الإحالات غير المترادفة لكل موقع غير مترادف

number of nonsynonymous substitutions per nonsynonymous site ورمزه K_A

يساوى

$$K_A = A_0 + [L_0 B_0 + L_2 B_2 / L_0 + L_2] \quad (٢٠)$$

وتباينه

$$V(K_A) = V(A_0) + [L_0^2 V(B_0) + L_2^2 V(B_2) / (L_0 + L_2)^2] - [2b_0 Q_0 (a_0 P_0 - c_0 (1 - Q_0))] /$$

$$L_0 + L_2 \quad (٢١)$$

٢.٢.٦. عدد إبدالات الأحماض الأمينية بين بروتينين.

يمكن حساب النسبة المشاهدة للأحماض الأمينية المختلفة بين تتابعين من

الأحماض الأمينية عند مقارنتهما على النحو التالى

$$p = n / L \quad (٢٢)$$

حيث n هي عدد الأحماض الأمينية بين التتابعين و L هي طول التتابعين المرصوين.
وباستخدام توزيع بواسون يمكن تحويل p إلى عدد إبدالات الأحماض الأمينية
لكل موقع على النحو التالي

$$d = -\ln(1 - p) \dots\dots\dots (٢٣)$$

وتباينه

$$V(d) = p / L(1 - p) \dots\dots\dots (٢٤)$$

٤.٦. رصن التتابعات في أزواج.

١.٤.٦. أساسيات رصن التتابعات.

تتضمن مقارنة تتابعين متماثلين تعيين مواضع الحذف والإدراج الذين يمكن أن يكونا قد وقعا في أي من التتابعين منذ تشعبهما عن سلفهما المشترك، وهو ما يطلق عليه عملية رصن التتابعات sequence alignment. وسوف نشرح هنا هذه العملية مستخدمين تتابعات الدنا مع العلم بأن عملية رصن تتابعات الأحماض الأمينية تخضع لنفس مبادئ وخطوات رصن تتابعات الدنا. بل في الحقيقة تعتبر النتائج المتحصل عليها من رصن تتابعات الأحماض الأمينية أجدر بالثقة من نتائج رصن تتابعات الدنا لسببين رئيسيين: (١) تتغير الأحماض الأمينية بصورة أقل من النيوكليوتيدات خلال التطور، (٢) هناك ٢٠ حمضا أمينيا بينما لا يتعدى عدد النيوكليوتيدات أكثر من ٤ ومن ثم فإن احتمال تماثل موقعين نتيجة للصدفة يكون أقل على مستوى الأحماض الأمينية عنه على مستوى النيوكليوتيدات

ويتكون رصن تتابع دنا من سلسلة من أزواج القواعد (قاعدة من كل تتابع)، ومن

ثم يكون ثمة ثلاثة أنواع من أزواج القواعد:

- **توافقات matches**، وفيها تظهر نفس النيوكليوتيدة في نفس الموقع في التتابعين ومن ثم نفترض عندئذ أن هذه النيوكليوتيدة لم تتغير منذ تشعب التتابعين من سلفهما المشترك.
- **عدم توافقات mismatches**، وفيها يكون في كل تتابع نيوكليوتيدة مختلفة في نفس الموقع، ومن ثم فإن إحلالا واحدا على الأقل قد وقع منذ تشعبهما من سلفهما المشترك.

• فجوات gaps، وتتكون من قاعدة موجودة في أحد التتابعين تقابلها قاعدة معدومة null base في نفس الموقع على التتابع الآخر، ويرمز للقواعد المعدومة بالرمز —، وتشير الفجوة إلى حدوث حذف في أحد التتابعين أو إدراج في الآخر. وعلى العموم لا ينبغي أن الرص في ذاته بأى منهما هو الذى قد وقع بالفعل.

وباعتبار وجود تتابعين من الدنا هما A وB طول كل منهما هو m و n على التوالي، فإذا رمزنا إلى عدد الأزواج المتوافقة بالرمز x وعدد الأزواج غير المتوافقة بالرمز y وعدد الفجوات بالرمز z فإن

$$n + m = 2(x + y) + z \quad (٢٥)$$

وعادة ما يتم التمييز بين الفجوات الطرفية terminal gaps والفجوات الداخلية internal gaps. فمثلاً يتم إقصاء الفجوات الطرفية من الحسابات عند مقارنة جزء من التتابع مع تتابع كامل، أو يتم إقصاء الفجوات الداخلية عند رص تتابع جينومى مع تتابع رنا رسول حيث أن الأول يكون محتويًا على انترونات وهو ما يفتقر إليه تتابع الرنا الرسول.

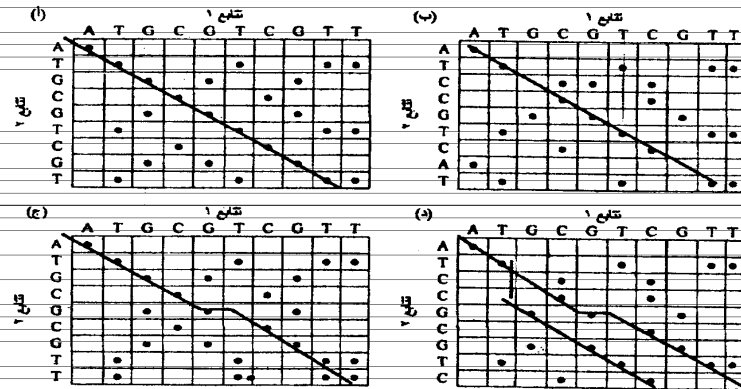
ومن حيث المبدأ ومن زاوية تطورية، فإن كل زوج في الرص يمثل تماثلاً موضعياً positional homology، أى أن عضوى الزوج منحدران من نيوكلويدية سلفية مشتركة. ومن ثم فإن أى خطأ في الرص يعنى بالضرورة غياب معنوية الرص حيث أننا نعجز حينئذ عن تقدير عدد الإحالات بصورة صحيحة.

فإذا علمنا أن الرص هو الخطوة الأولى في العديد من الدراسات الجينومية وأن أى خطأ في الرص سوف يتضاعف في العمليات الحسابية اللاحقة، فإنه يتحتم علينا أن نقوم بعملية الرص بمنتهى الدقة. وعلى المرء أن يتخلص من كل الأجزاء المثيرة للشك من الرصيص قبل أن يبدأ في التحليل، حتى وإن أدى ذلك إلى نقص طول الرص بصورة ملحوظة وما يصحبه من ارتفاع في قيمة خطأ المعاينة عند تقدير عدد الإحالات النيوكليوتيدية بين التتابعين.

٢.٣. رسم التتابعات باستخدام طريقة مصفوفة النقاط.

لعل طريقة مصفوفة النقاط dot matrix method التي وضعها Gibbs and McIntyre (1970) هي أشهر وأسهل طرق رسم التتابعات، وفيها تتم كتابة التتابعين قيد الدراسة كعناوين رأسية وأفقية لأعمدة وصفوف مصفوفة ذات بعدين (شكل ٨-٦)، ثم توضع نقطة على رسم مصفوفة النقاط dot matrix plot عند موضع تماثل النيوكليوتيدات في التتابعين. ومن ثم فإن نقطة عند (x,y) تعني أن نيوكليوتيدة عند الموقع x في التتابع الأول هي نفسها النيوكليوتيدة عند الموقع y في التتابع الثاني.

ويعرف الرصيص alignment بأنه طريق في المصفوفة يبدأ من أعلى العناصر إلى اليسار وينتهي عند أسفل العناصر إلى اليمين. ومن ثم فإن هناك أربع خطوات لهذا الطريق : (١) خطوة قطرية عبر نقطة تعني توافقاً، (٢) خطوة



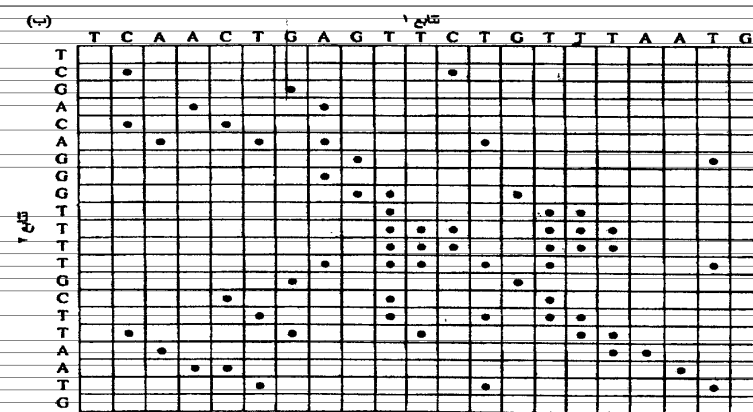
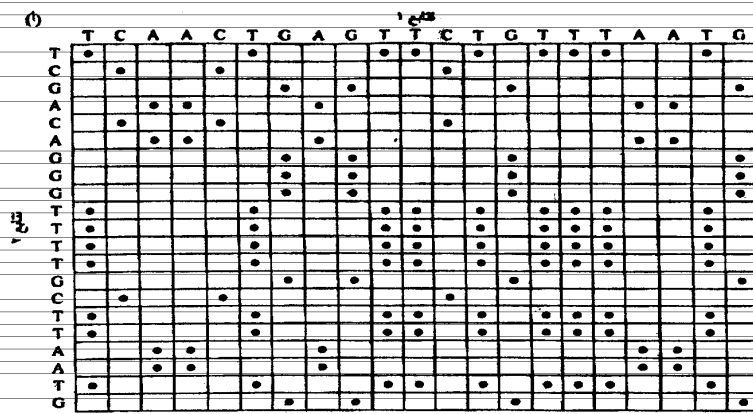
شكل (٨-٦) - مصفوفات نقاط لرص تتابعين نيوكليوتيديين. (أ) التتابعان متماثلان في الجزء

المخصوص. (ب) يختلف التتابعان عن بعضهما البعض ولكن لا يحتوى الجزء المخصوص على فجوات. (ج) يحتوى الرصيص على فجوة، ولكن فيما عدا ذلك يتماثل التتابعان. (د) يظهر رصيصان يحتوى كل منهما على فجوات وعدم توافقات مما يدفع إلى استخدام بعض الحسابات للاختيار بينهما.

قطرية عبر عنصر خال في المصفوفة تعنى عدم توافق، (٢) خطوة افقية تعنى وجود نيوكليوتيدة معدومة (فجوة) في التتابع على رأس المصفوفة، (٤) خطوة رأسية تعنى وجود نيوكليوتيدة معدومة (فجوة) في التتابع على جانب المصفوفة.

فإذا كان التتابعان تامي التماثل (أو إذا تمت مقارنة تتابع بنفسه)، فسيتم وضع نقاط في كل العناصر القطرية في المصفوفة (شكل ١٨-٦). أما إذا اختلف التتابعان عن بعضهما البعض فقط عن طريق إحلالات فسيتم وضع نقاط على معظم العناصر القطرية في المصفوفة (شكل ١٨-٦ ب). أما إذا حدث إدراج في أحد التتابعين بلا أي إحلال فسوف توجد منطقة داخل المصفوفة يتم فيها نقل قطر الرصيص إما رأسيا أو أفقيا (شكل ١٨-٦ ج). وفي كل من الحالات السابقة الثلاث، كان الرصيص يشير إلى نفسه، أما إذا اختلف التتابعان نتيجة لوجود فجوات وإحلالات بينهما، فإنه يصعب تعيين موقع الفجوات ومن ثم يتحتم علينا الاختيار من بين أكثر من رصيص (شكل ١٨-٦ د)، وفي مثل هذه الحالات يفضل استخدام طرق أكثر ثقة من طريقة مصفوفة النقاط.

وبلاحظ في الشكل (١٩-٦) أنه في معظم الحالات تكون مصفوفة النقاط شديدة التشوش أي يتم شغل عناصر أخرى في المصفوفة غير تلك المثلة للرصيص الحقيقي بنقاط تغطي على الرصيص، حيث أنه عند مقارنة تتابعين من الدنا، فإن حوالي ٢٥٪ من عناصر المصفوفة يتم شغلها بالصلفة وحدها. وثمة مقياسان لتحديد عدد التوافقات الزائفة ومن ثم درجة وضوح رسم مصفوفة النقاط وهما حجم النافذة window size والشرطية stringency. وعليه فإنه بدلا من استخدام مواقع نيوكليوتيدية مفردة، فإنه يمكن مقارنة التتابعات باستخدام نوافذ متداخلة (منزلة) ذات أطوال ثابتة، وكل مقارنة داخل المصفوفة لا بد أن تحقق قيمة حدية دنيا معينة لمجموعة النافذة ككل (الشرط) حتى يتسنى اعتبارها توافقا. فمثلا في الشكل (١٩-٦ ب) نجد أننا قد استخدمنا حجم نافذة طوله ثلاثة أزواج من القواعد، ولا توضع نقطة في المصفوفة إلا بشرط أن تتوافق قاعدتان على الأقل من الثلاث بين التتابعين، وبهذا تصبح المصفوفة أقل تشوشا.



تتابع ١

(ج)

	S	T	E	F	C	L	M
S	•						
T		•					
G							
F				•			
C					•		
L						•	
M							•

شكل (٩-٦) - (أ) مصفوفة نقاط لتتابعين نيوكليوتيديين، الرصيص مقطى بالعديد من النقاط الزائفة في المصفوفة. (ب) مصفوفة النقاط لتتابعي النيوكليوتيدات الموجودين في (أ)، ولكن بعد استخدام حجم نافذة من ثلاث نيوكليوتيدات ووضع حد شرطية باثنين من ثلاثة توافقات لوضع نقطة في العنصر المناسب، وتساعد عملية الفربلة هذه على حذف عدد كبير من النقاط الزائفة ومن ثم يظهر الرصيص بصورة أوضح على خلفية أقل تشوشاً. (ج) مصفوفة نقاط لتتابعين من الأحماض الأمينية تم الحصول عليهما بترجمة تتابعي النيوكليوتيدات الموجودين في (أ)، ويتضح كيف أن الرصيص لم يعد غامضاً بالمرّة.

وفي حالة التتابعات المشفرة للبروتينات، فإنه يفضل مقارنة تتابعات الأحماض الأمينية بدلا من مقارنة تتابعات الدنا، حيث أن زيادة الصفات (الحروف) من ٤ إلى ٢٠ ونقص طول التتابع من L إلى L/3 سوف يخفض بشدة من عدد النقاط الزائفة (شكل ٩ج).

٢.٤.٦. حسابات رصن التتابعات.

لقد رأينا كيف أن استخدام طريقة مصفوفة النقاط لرص تتابعات ذات إحالات وفجوات فإنها تعطي أكثر من رصيص، ولكن أيا منها هو الرصيص الأمثل optimal alignment، أي أفضل رصيص محتمل بين تتابعين، الذي يمكن الاعتماد عليه. والرصيص الأمثل هو ذلك المحتوى على الحد الأدنى من عدم التوافقات والفجوات تبعا لمعايير معينة. بيد أنه وللأسف، يؤدي خفض عدد عدم التوافقات إلى زيادة عدد الفجوات والعكس بالعكس.

فمثلا إذا كان لدينا التتابعان التاليان A و B :

$L_A = 11$	TCAGACGATTG	A:
$L_B = 9$	TCGAGACTG	B:

فإنه يمكننا خفض عدد عدم التوافقات إلى صفر كما يلي :

TCAG-ACG-ATTG

(I)

TC-GGA-GC-T-G

ونحصل هنا على ٦ فجوات. وعلى النقيض من ذلك يمكننا خفض عدد الفجوات إلى فجوة واحدة لها أقصر طول مسموح به وهو هنا نيوكليوتيدتان حيث:
 $|L_A - L_B| = 2$ ، مما يترتب عليه زيادة في عدد عدم التوافقات :

TCAGACGATTG

* * * * * (II)

TCGGAGCTG--

وفي هذه الحالة نحصل على فجوة واحدة (طولها نيوكليوتيدتان)، بيد أن عدد عدم التوافقات (المشار إليها بنجمة) يرتفع إلى ٥ (١ من نوع الـ transition و٤ من نوع الـ transversion).

وكبديل، يمكننا اختيار رصيص لا يخفض لا من عدد الفجوات ولا من عدد عدم التوافقات، مثل :

TCAG-ACGATTG

* * (III)

TC-GGA-GCTG-

وفي هذه الحالة يكون عدد عدم التوافقات ٢ (كلاهما من نوع الـ transversion) وعدد الفجوات ٤.

فمن إذن من بين الرصاص الثلاثة هو الأمثل ؟ ولأن مقارنة الفجوات مع عدم التوافقات تشبه مقارنة التفاح بالبرتقال أي لا معنى لها، ومن ثم فإنه يتحتم علينا إيجاد قاسم مشترك يمكن بواسطته مقارنة الفجوات بعدم التوافقات، وهو ما يطلق عليه غرامة الفجوة أو تكلفة الفجوة gap penalty or gap cost. وغرامة الفجوة هي عامل (أو عدة عوامل) تضرب بها قيم الفجوة (أعداد الفجوات وأطوالها) حتى تتساوى قيمة الفجوات مع

قيمة عدم التوافقات. وتبنى غرامة الفجوات على تقديرنا لتكرارات الأنواع المختلفة من الإدراج والحذف التي وقعت خلال التطور مقارنة بتكرار الإحلالات النيوكليوتيدية. وعليه فإنه لا بد لنا كذلك من تقدير غرامة عدم التوافقات mismatch penalties الذى هو تقدير لتكرار الإحلالات.

ولأى رصيص، يمكننا حساب مسافة أو مؤشر عدم تشابه distance or dissimilarity index بين التتابعين فى الرص ويرمز له بالرمز D

$$D = \sum m_i y_i + \sum w_k z_k \dots\dots\dots (٢٦)$$

حيث y_i عدد عدم التوافقات من النوع أ، و m_i غرامة عدم التوافق من النوع أ، و z_k عدد الفجوات ذات الطول k، و w_k رقم موجب يمثل غرامة الفجوة ذات الطول k.

وبالمثل، فإن التشابه بين تتابعين فى رص يمكن قياسه بمؤشر تشابه similarity index يرمز له بالرمز S، وهو يساوى

$$S = x - \sum w_k z_k \dots\dots\dots (٢٧)$$

حيث x عدد التوافقات.

وعادة ما يتم افتراض احتواء غرامة الفجوات على مكونين هما غرامة فتح الفجوة gap-opening penalty وغرامة إطالة الفجوة gap-extension penalty، وتعتمد الأخيرة على معلومية مسبقة بتكرار أحداث الحذف والإدراج ذات أطوال محددة بالنسبة لأحداث الإحلالات النيوكليوتيدية الأخرى. وفى نظام الغرامة الثابتة للفجوات fixed gap penalty system لا يتم حساب أى غرامة لإطالة الفجوات، أما فى نظام

الغرامة المتصلة أو الخطية للفجوات affine or linear gap penalty system تحسب تكلفة إطالة الفجوة بضرب طول الفجوة مطروح منها ١ فى ثابت يمثل غرامة إطالة الفجوة ب١. فمثلاً، بالنسبة لفجوة طولها ١ تتضمن تكلفة الفجوة غرامة فتح الفجوة فقط، أما بالنسبة لفجوة طولها ٣ فإن تكلفة الفجوة تتضمن غرامة فتح الفجوة بالإضافة إلى غرامة إطالة الفجوة مضروبة فى ٢. وحتى لا تزيد غرامة إطالة الفجوة بصورة كبيرة فى الفجوات الطويلة، فقد اقترح بعض العلماء مثل (Gu and Li (1995 نظام الغرامة

اللوغاريتمية للفجوات logarithmic gap penalty system، وفيه تزداد غرامة الفجوة بصورة أبطأ مع الزيادة في طول الفجوة.

وبالنسبة للرصاص الثلاثة السابقة، فإننا إذا استخدمنا غرامة عدم توافق تساوى ١ وغرامة فتح فجوة تساوى ٢ وغرامة إغلاق فجوة تساوى ٦، يصبح مؤشر عدم التشابه للرصاص ١ يساوى $12 = 6(1-1) + (6 \times 2) + (0 \times 1)$ ، وبالمثل، فإن قيمة D للرصاصين ١١ و ١٢ تساوى 10 و 12 على التوالي. وعليه فإنه من بين الرصاص الثلاثة هذه يعتبر الرصاص ١١ هو الرصاص الأمثل (أقل D).

إن الفرض من أى حساب للرص alignment algorithm هو اختيار الرصاص ذى الـ D الأقل (أو الـ S الأعلى) من بين جميع الرصاص المحتملة الأخرى، والتي عادة ما تبلغ أعدادها أرقاما فلكية. فمثلا إذا تمت مقارنة تتابعين طول كل منهما ٢٠٠ فضالة فإننا نحصل على 10^{40} رص محتمل وذلك بالسماح بأى عدد وأى طول من الفجوات. ومن حسن الحظ فتمت حسابات كمبيوترية للعثور على الرصاص الأمثل لعل أشهرها حساب نيدلمان وونش (Needleman-Wunsch algorithm (1970) والذي يستخدم تقنية كمبيوترية عامة تدعى البرمجة الديناميكية dynamic programming وهى تقنية تستخدم حين يمكن تقسيم بحث كبير إلى عدد من الخطوات من المراحل الصغيرة بحيث يكون (١) حل مرحلة البحث المبدئية عاديا، (٢) كل حل جزئى للمراحل المتأخرة من البحث يمكن إيجاده بالرجوع إلى عدد قليل فقط من الحلول فى المراحل الأبر، (٣) تحتوى المرحلة الأخيرة على الحل الكلى. ومن ثم يمكن إجراء البرمجة الديناميكية على فضايا الرص لأن مؤشرات التشابه تخضع للقاعدة التالية :

$$S_{1 \rightarrow x, 1 \rightarrow y} = \max S_{1 \rightarrow x-1, 1 \rightarrow y-1} + S_{x,y} \quad (٢٨)$$

حيث $S_{1 \rightarrow x, 1 \rightarrow y}$ هى مؤشر التشابه بين التتابعين حتى الفضالة X فى التتابع الأول والفضالة Y فى التتابع الثانى، و $S_{1 \rightarrow x-1, 1 \rightarrow y-1}$ هو مؤشر التشابه لأفضل تتابع حتى الفضالة $x-1$ فى التتابع الأول والفضالة $y-1$ فى التتابع الثانى، و $S_{x,y}$ هو نتيجة التشابه لرص الفضالتين X وY.

وهو ما يصدق بحذافيره كذلك على مؤشر عدم التشابه بحيث :

$$D_{1 \rightarrow x, 1 \rightarrow y} = \max D_{1 \rightarrow x-1, 1 \rightarrow y-1} + D_{x, y} \quad (٣٨)$$

ويتم حساب الرص على مرحلتين :

أولاً، يتم رص التتابعين بطريقة مصفوفة النقاط، ولكل عنصر في المصفوفة، وليكن x

و y ، يحسب مؤشر التشابه $S_{1 \rightarrow x, 1 \rightarrow y}$. وفي نفس الوقت، يتم تخزين أفضل نتيجة للرص

في الصف أو العمود السابق، وتسمى القيمة المخزنة بالمؤشرة pointer. وتمثل العلاقة

بين قيم $S_{1 \rightarrow x, 1 \rightarrow y}$ الجديدة والمؤشرة بسهم.

ثانياً، يتم إنتاج الرصيصة بدءاً بأعلى نتيجة تشابه سواء في العمود الواقع أقصى اليمين أو

في الصف السفلى ثم يتم تعقب أفضل مؤشرة من اليمين إلى اليسار، ويطلق على هذه

العملية القيافة traceback، كذلك يسمى رسم المؤشرات على خط القيافة برسم

الطريق path graph لأنه يعرف الطرق عبر المصفوفة التي توافق الرصيصة الأمثل.

ويظهر الشكل (١٠-٦) مثالا بسيطا لبرمجة ديناميكية لرص التتابعين ATGCG و

ATCGCG، ولتبسيط الأمر فإننا نفرض أن $w_k = 0$ (أي لا غرامة للفجوات)، ومن ثم فإن

الرصيصة الأمثل هو ذلك الرصيصة ذو العدد الأكبر من التوافقات. وتملأ المصفوفة من

اليسار إلى اليمين ومن أعلى إلى أسفل، ومن ثم نجد توافقا واحدا في الصف الأول يحصل

على النتيجة ١ وأربعة عدم توافقات كل منها يحصل على النتيجة صفر، أما في الصف

الثاني فإننا نبدأ بمزاوجة T من التتابع الأول ب-T من التتابع الثاني وهو توافق ومن ثم

فإننا نضيف ١ إلى أعلى نتيجة سابقة (وهي ١ في الصف الأول) ونكتب المجموع عندها مع

رسم سهم من هذا العنصر إلى مؤشرته، وهكذا دواليك إلى أن يتم الانتهاء من المصفوفة كلها

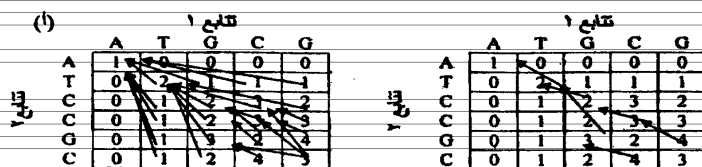
(شكل ١٠-٦)، مع ملاحظة أنه يمكن لعنصر واحد أن يكون له أكثر من مؤشرة أو أن تشير

أسهم أكثر من عنصر إلى مؤشرة واحدة. وفي المرحلة الثانية، نبدأ من أعلى نتيجتين في

آخر صف وآخر عمود من المصفوفة (مظللتين في الشكل ١٠-٦ أ ب) ونبدأ في اقتفاء

الرصيصة بالبحث عن المؤشرة (أو المؤشرات) السابقة ذات القيمة الأعلى، حتى نحصل على

طريق أو طرق تبرز الرصائص المثلى.



شكل (١٠-٦) -- حساب البرمجة الديناميكية لرص على مرحلتين (أ وب)، ويلاحظ الحصول على رصيصين أمثلين في (ب).

٥.٦. رصن التتابعات المتعددة.

١.٥.٦. حسابيات وبرامج الرصائص المتعددة.

يمكن اعتبار رصن التتابعات المتعددة multiple sequence alignment

امتدادا لرصن التتابعات في أزواج، بيد أن حساباته تزداد تعقيدا وبصورة أسية كلما زاد عدد التتابعات قيد الدراسة، ومن ثم يفسح من غير العمل إجراء بحث شامل عن الرصيص الأمثل في حالة رصن التتابعات المتعددة. ولقد تم نشر عدد من الطرق المساعدة على ذلك، كما تم وضع عدد من برامج الكمبيوتر المتاحة على الانترنت التي تقوم بالعثور على الرصيص الأمثل لتتابعات متعددة (انظر القائمة التالية):

1- *Baylor's search launcher for Biologists*, Multiple Alignment section
ClustalW, MAP, PIMA 1.4, MSA 2.1, and BLOCK MAKER alignments
<http://dot.imgen.bcm.tmc.edu:9331/multi-align/multi-align.html>

2- Clustalw + Multalin Multiple Alignment at IBCP

Calculate, optimize and score multiple sequence alignments with ProbModel, MSA, ClustalW and heuristics. The scoring is a probabilistic measure.
<http://cbrg.inf.ethz.ch/MultAlign.html>

3- STRAP: Interactive program for generating and analyzing multiple sequence alignments. Multiple structure alignments, integrated 3D-viewer, 3D-superposition of protein backbones, mapping of mutations

onto 3D-models, translation of nucleotide sequences to amino acid sequences, sequence dotplots.

<http://www.charite.de/bioinf/strap/>

4- MUSCA: An Algorithm for Constrained Alignment of Multiple Data Sequences, and related tools.

<http://cbcsrv.watson.ibm.com/Tmsa.html>

5- JEvTrace - algorithms and multivalent graphical browser for combined alignment, phylogeny, and structure analysis.

<http://www.cmpharm.ucsf.edu/~marcini/JEvTrace/>

6- MAVID - multiple alignment program for large genomic regions

<http://baboon.math.berkeley.edu/mavid/>

7- MUSEQAL optimal multiple alignment by iteratively improving a given set of pre-aligned sequences.

<http://godzilla.dcrf.nih.gov/~yap/museqal2.html>

8- DCA: Divide and Conquer Multiple Sequence Alignment

Close-to-optimal, fast simultaneous sum-of-pairs alignment of up to twenty protein, DNA, or RNA sequences

<http://bibiserv.techfak.uni-bielefeld.de/dca/>

9- DIALIGN - fragment-based alignment program following the method described by Morgenstern *et al.* in *PNAS* 93 (22), pp. 12098-12103.

<http://cartan.gmd.de/ToPLign.html>

10- ToPLign: Toolbox for Protein Alignment

Computing, analysis and visualization of pairwise, multiple, threading, and parametric alignments.

<http://www.ibc.wustl.edu/msa.html>

11-MSA, (Close-to-) Optimal Alignments using the Carrillo-Lipman bound

<http://www.toulouse.inra.fr/multalin.html>

تشابها هما تتابعى ميوجلوبين الحوت (الثالث من أسفل - رقم ٥) ولجهيموجلوبين نبات الترمس البقولي (التتابع الأخير- رقم ٧)، حيث لا يتعدى التماثل بينهما أكثر من ١٠٪ ومع ذلك فإن التشابه بينهما يظل ظاهرا عند رصهما.

وتجدر الإشارة إلى أنه للحصول على مثل هذا الرصيص فلا بد أولا من التأكد من صحة ودقة الحصول على البيانات الجزيئية (التتابعات) حيث أنه ما من برنامج يقدر على إنتاج رصيص صحيح من بيانات غير صحيحة، وهى مهمة يضطلع بها باحثو البيولوجيا الجزيئية قبل استخدام برنامج الكمبيوتر، حيث عليهم توخى الحذر والنظر إلى التتابع بفرض حذف الأجزاء المستولة عن تشويه النتائج. ويمكن تحقيق هذا باتباع مدخل iterative، فيقوم الباحث بعمل عدة رصائص اختبارية ويختبر كلا من الشجرة الفيلوجينية المتحصل عليها علاوة على الرصائص البينية.

وأخيرا، إذا توفرت للباحث معلومات إضافية مثل معرفة فضالات موقع نشط محدد فإن مثل هذه المعلومات من شأنها أن تساعد على إلقاء الضوء على بعض أخطاء الرص أو على التحرير اليدوى للرصيص. وهو ما يجعل الباحث فى حاجة إلى استخدام نوعين من البرامج أحدهما ليولد رصيصا مناسباً والآخر لتحريره. وعادة ما يتم اختيار التتابعات من قواعد البيانات الموجودة على الإنترنت مثل (Bairoch and Apweiler, 1996) SwissProt باستخدام برامج البحث على أساس التشابه مثل (FASTA (Pearson and Lipman, 1988 أو BLAST (Altschul *et al.*, 1997). وكذلك يمكن إجراء البحث باستخدام مفاتيح keywords أو أدوات البحث داخل قواعد البيانات مثل SRS (Etzold *et al.*, 1996) أو كليهما. والطريقة الأولى دوما ما تنتج تتابعات شبيهة بالتتابع قيد الدراسة (والذى قد تم الحصول عليه معمليا أو من أحد الأدبيات) حتى وأن كان هذا التشابه لا يعكس أى علاقة تطورية وإنما هو وليد تركيب الأحماض الأمينية الموجودة فيهما، ومن ثم فعلى الباحث إعمال عقله ومعرفته بالخلفية البيولوجية للتابع قيد الدراسة حتى يتسنى له اختيار التتابع الصحيح للرص.

وفى الحالة الأولى، قد يساعد رمز التتابعات المتعددة الباحث فى معرفة ما إذا كان أى توافق هو توافق صحيح وليس وليداً للصدفة، فكما يظهر فى الشكل (١١-٦) فإننا نجد أنه ثمة عدة فضالات محفوظة بين كل الجلوبيينات فى الرصيصة (مثل هستيديني ربط الهيم)، حيث تبدو هذه الفضالات كحزم رأسية محفوظة وهو ما لا يمكن اكتشافه عند رمز تتابعين فقط، فنحن مثلاً لا نستطيع البت بوجود تشابه من عدمه بين ميوجلوبيين الحوت ولجهيموجلوبيين التمس عند رصنهما وحدهما، بيد أن هذا التشابه يظهر جلياً عند استخدام رمز التتابعات المتعددة. ولكن لا بد من الانتباه إلى إذا ما كانت التتابعات المتحصل عليها باستخدام التشابه تحتوى على شظايا قصيرة كمجالات غير تامة أو على مجالات متعددة كالفيرونكتين والـ SH3 حيث أن هذا سوف يسبب مشاكل بالنسبة لصحة الرصيصة ومن ثم يفضل رمز المجالات كل على حدة إلا إذا احتوت كل التتابعات على نفس المجالات بنفس الترتيب.

أما فى الحالة الثانية عند استخدام مفاتيح بحث بدلا من تشابه التتابعات، فلا بد من معرفة أن وجود تتابعين لهما نفس مفتاح البحث لا يعنى بالضرورة وجود علاقة تطورية بينهما. فعلى سبيل المثال قد تقوم بروتينات مختلفة بنفس الوظيفة الإنزيمية فى مجموعات الكائنات الحية المختلفة، أو يكون مصحح قاعدة البيانات أو العالم نفسه قد قاما بخطأ عند تسمية التتابع.

٢.٥.٦. استخدام برنامج CLUSTAL W فى رمز التتابعات.

تكمن قيمة برنامج CLUSTAL W فى رمز التتابعات المتعددة للنسب والبروتينات فى أنه يعطى رصائص ذات دلالات بيولوجية للتتابعات العديدة المتشعبة من بعضها البعض، حيث يقوم بحساب أفضل توافق للتتابعات المختارة ثم رمز هذه التتابعات بصورة يمكن معها مشاهدة التشابهات والاختلافات بينها وذلك بتعيين المناطق المحفوظة فيها، وهو ما يعد خطوة أساسية عند تصميم التجارب لاختبار وتحوير وظيفة بروتين ما أو عند التنبؤ بتركيب البروتينات ووظيفتها وكذلك عند تعيين أعضاء جدد فى العائلات البروتينية. كذلك يقوم البرنامج بتوضيح العلاقات التطورية بين التتابعات فى صورة شجرة فيلوجينية phylogram أو طائفية cladogram.

وبرنامج CLUSTAL W متاح على عدة مواقع على الانترنت منها <http://www.ebi.ac.uk/clustalw/>. وثمة طريقتان لاستخدامه على هذا الموقع، الأولى تفاعلية interactive، حيث ينتظر المستخدم ظهور النتائج على نافذة الاستعراض browser window، أما الطريقة الثانية فهي عن طريق البريد الإلكتروني email حيث يتم إرسال النتائج على البريد الإلكتروني للمستخدم، وتعتبر الطريقة الثانية هي المفضلة عند رسم عدد كبير من التتابعات. كذلك يمكن الحصول على CLUSTAL W بطرق أخرى.

ويمكن إدخال التتابعات إلى البرنامج بأحد سبعة تنسيقات وإن كانت Fasta وALN/ClustaW هما أكثرها استخداما، أما نتائج الرسم فيحصل عليها بتنسيق ALN/ClustaW أى بالامتداد .aln، وهو النسق الذى يتيح تلوين الرصيص فى حالة البروتينات تبعا لخواصها الكيموفيزيائية على النحو الموضح فى الجدول التالى:

AVFPMILW	RED	Small (small+ hydrophobic (incl.aromatic -Y))
DE	BLUE	Acidic
RHK	MAGENTA	Basic
STYHCNGQ	GREEN	Hydroxyl + Amine + Basic - Q
Others	Gray	

ويستخدم هذا الجدول عند كتابة الرصيص، وكذلك يدل الرمز "*" على

التوافق، والرمز ":" على حدوث إحلال محفوظ، والرمز "." على إحلال نصف المحفوظ.

٧. الجينومكس المقارن و الفيلوجينيا

Comparative Genomics & Phylogenetics.

إعداد: سناء رياض و أحمد المتينى

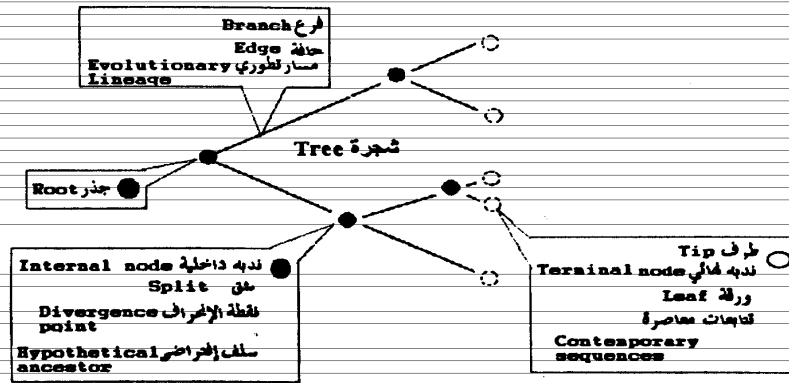
إن الدراسات المبنية على تتبع علاقات القرابة بين الكائنات (phylogenetic) اعتمادا على تحليل الـ DNA أو البروتينات، اكتسبت أهمية في السنوات القليلة الماضية بفرض دراسة العلاقات التطورية للكائنات من أول البكتيريا حتى الإنسان. وحيث أن معدل التغير لهذه التتابعات (على مستوى الـ DNA أو البروتين) عال وكبير بشكل ملحوظ (Nei, 1996)، فإننا نستطيع في حقيقة الحال دراسة العلاقات التطورية على كافة المستويات التقسيمية للكائنات (kingdom, phyla, classes, families, genera, species, and populations). ولتعدد هذه الوحدات التقسيمية، لجأ العلماء لتحديدها (على اختلاف أشكالها) بمصطلح الوحدة التقسيمية (taxon)، والتي قد تكون كائنا أو جينا أو حتى تتابع من الـ DNA أو من البروتين، والتتابعات الجزيئية الأخيرة ستمثل أساس تناولنا لهذا الموضوع حيث أنها أكثر ارتباطا بمفهوم المعلوماتية الحيوية. كذلك فإن مثل تلك الدراسات يمكن أن تلقى الضوء على تطور العائلات الجينية (gene families) ويمكنها أيضا أن تلقى ببعض الضوء على ميكانيكيات الحفاظ على تعدد الأشكال المظهرية (polymorphic genes) على المستوى الجزيئي.

١.٧. الشجرة التطورية Evolutionary tree.

إن الشجرة التطورية أو شجرة القرابة (phylogenetic tree) تعرف أيضا باسم dendrogram ويعنى العلاقة بين مجموعة من الأفراد المختلفة وإذا كانت تلك العلاقات فيولوجينية تسمى العلاقة cladogram. وقد بدأت تلك الدراسات منذ ٤٠ سنة مع بدايات علم التقسيم الرقمى (numerical taxonomy) المعتمد على الصفات المظهرية التقليدية

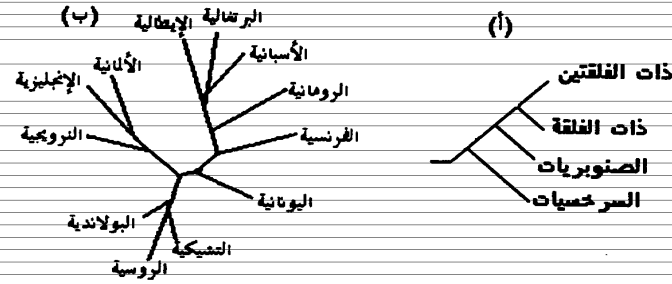
وعلى دراسة التكرارات الجينية في العشائر، وبعض من برامج التحليل التي وضعت في السابق لهذه الأغراض مازالت تستعمل حتى اليوم لتحليل البيانات الجزيئية، وأن كان لدينا الآن برامج أكثر حداثة وكفاءة لهذا الغرض.

والتقنية العملية لتصميم شجرة القرابة، تتبع الأساليب الإحصائية لعمل شجرة اعتيادية ممثلة للشجرة الحقيقية، ولهذا فلا بد من معرفة طوبوغرافية (topology) الشجرة أو نظام تشعبها أولاً ثم يتبع ذلك حساب مدى أو طول هذا التشعب. والشجرة عبارة عن رسم ثنائي الأبعاد من فروع (branches) تمثل درجات القرابة بين الـ taxa، والتي قد تكون مرتبطة بمنشأ مشترك واحد يسمى عادة بالجنر root تخرج منه الفروع branches الممثلة للمسارات التطورية ويوجد عليها ندى إما داخلية internal nodes وإما طرفية نهائية terminal nodes وهي تمثل البراعم الجانبية والبراعم الطرفية للشجرة والتي يمكن اعتبارها الأوراق leaves. وشكل (١-٧) يمثل شجرة قياسية مبين عليها المصطلحات والمرادفات المستخدمة.



شكل (١-٧) : رسم توضيحي لشجرة فيلوجينية قياسية.

وقد تكون الشجرة ذات جذر (منشأ) rooted tree كما هو مبين بشكل (١٢-٧) أو
عديمة الجذور un-rooted tree وهو ما يعرف أيضا بأسم الشجرة النجمية star tree كما
هو مبين بشكل (٢-٧).



شكل (٢-٧) : (أ) مثال للشجرة ذات الجذور للعلاقة بين النباتات الراقية، (ب) مثال للشجرة
النجمية للعلاقة بين اللغات الأوروبية.

٢.٧ طرق تعيين الشجرات Tree determination methods

إن بناء طوبوغرافية الشجرة، يتطلب قدرا كافيا من البيانات الجزيئية (افترض
أنها من عشر مصادر) لذا فإن عدد تفرعات الشجرة المحتملة قد يصل إلى المليون وعليك أن
تختار منها الأنسب. ولقد اقترحت لذلك عشرات الطرق الإحصائية المعتمدة على
الكمبيوتر لتحقيق ذلك ويمكن إجمالها تحت ثلاث أقسام رئيسية هي :

- ١- الطرق المقتصدة parsimony methods.
- ٢- طرق حساب المسافة (البعد) distance methods.
- ٣- طرق حساب أرجحية التشابه likelihood methods.

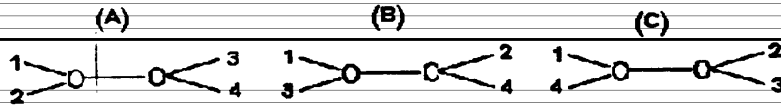
١.٢.٧ الطرق المقتصدة Parsimony Methods

إن الطرق المقتصدة لعمل الشجرات التطورية بين مجموعة من التتابعات تعتمد
أساسا على إمكانية اختيار العدد الأدنى من التغيرات في التتابعات بين الـ taxa تحت
الدراسة. ولعمل الشجرة باتباع هذه الطريقة لابد أولا من عمل الصف المتعدد

(multiple alignment) للعينات تحت الدراسة، والشجرة التي يمكن أن تفسر الاختلافات الموجودة عند أقل عدد ممكن من التبادلات هي التي يتم اختيارها. وهذه الطرق مطولة وتحتاج لجهود ووقت كبير، وعادة تنجح في دراسة التتابعات التي بينها درجات عالية من القرابة. ولتوضيح هذه الطرق دعنا نناقش المثال التالي، فالتتابعات التالية تمثل أربعة عينات مختلفة وتتابعاتها هي:

1	AAGAGTGCA
2	AGCCGTGCG
3	AGATATCCA
4	AGAGATCCG

هذه التتابعات الأربع يمكن أن نكون فيما بينها ثلاث شجرات نجمية (ليس لها جنس) كما هو موضح بشكل (٢-٧).



شكل (٢-٧) ، الشجرات النجمية الثلاث التي يمكن أن تمثل العلاقات بين التتابعات الأربع المبينة.

و يجب أن نلفت النظر إلى ملاحظة هامة وهي أن التتابعات السابقة يوجد بها صفوف مفيدة معلوماتيا informative وأخرى غير معلوماتية un-informative. الصف المعلوماتي هو الموضع الذي يلاحظ به اثنين من القواعد المتماثلة على الأقل، وعلى ذلك فجميع المواضع السابقة معلوماتية ما عدا الموضع الرابع حيث لا توجد قواعد متماثلة على الإطلاق. وعلى البرنامج الحسابي (برنامج الكمبيوتر) أن يميز الشجرة الأكثر كفاءة منها حيث يوجد أقل عدد ممكن من التغيرات بين التتابعات الأربع. وعلى ذلك يمكن القول بأن الشجرة الأولى (A) هي الأكثر كفاءة لتمثيل العلاقة بينهم حيث عدد التغيرات سيكون ٢ و ١ فقط عند مواقع التلاقي على الشجرة، بينما النماذج الأخرى ستعطي عدد أكبر من التغيرات.

وهذه الطرق تعتبر من أقدم الطرق المستخدمة لتقدير شجرات القرابة لكن كفاءتها محدودة خصوصا عند التعامل مع عدد كبير من الوحدات التقسيمية مستخدمين تناوبات كبيرة للصف والمقارنة.

لكن يجب التنويه بأن البرامج الحسابية أو برامج الكمبيوتر المختصة بتقدير وتصميم الشجرات التقسيمية التطورية تعتمد في بعض منها على نظرية من نظريات علم وراثية العشائر population genetics ألا وهي نظرية "الزمن الجزيئي" molecular clock التي ظهرت في أواخر الستينات من القرن الماضي. فمن دراسات تغيرات (طفرات الإستبدال) الأحماض الأمينية في البروتينات، إتضح أن معدل الإستبدالات ثابت تقريبا مما يعنى أن التفرعات على الشجرة إذا ما تساوت في أطوالها branch length يعنى التساوى في عدد الإستبدالات المحتملة وعليه يمكن تحديد زمن التفرع على الشجرة أو زمن الفصل بين الوحدات التقسيمية. وقد لاقى هذه النظرية قبولا كبيرا بين علماء وراثية العشائر والتطور الذين يؤمنون بحيادية الإختلافات الجزيئية neutrality وعلى رأسهم Kimura، وإن لاقى عدم قبول من بعض العلماء المشككين من مبدأ الحيادية هذا. وعليه فهناك برامج حسابية algorithm تأخذ في الحسبان نظرية الزمن الجزيئي وغيرها تتجاهلها.

٧.٢.٢. طرق تقدير البعد Distance methods.

هناك العديد من الطرق التي تعتمد على تقدير البعد الوراثي بين ال taxa تحت الدراسة، وسنحاول هنا تناول أهمها بالشرح والتوضيح. فعمل الشجرات إعتقادا على البعد، تقدر العلاقات بين العينات إما بحساب عدد الإستبدالات على الفرع أو بحساب عدد التغيرات الواجب إجرائها لجعل أى تنابعين متناظرين مرة أخرى. ونجاح هذه الطرق يعتمد على قدرتها في تبسيط الإختلافات وتحويلها إلى صورة قيم مضافة additive. ولتوضيح الفكرة الأساسية من تلك الطرق دعنا ندرس التناوبات التالية بين أربع من العينات لعمل شجرة تربطهم:

1. ACGCGTTGGGCGATGGCAAC
2. ACGCGTTGGGCGACGGTAAT
3. ACGCATTGAATGATGATAAT
4. ACACATTGAGTGATAATAAT

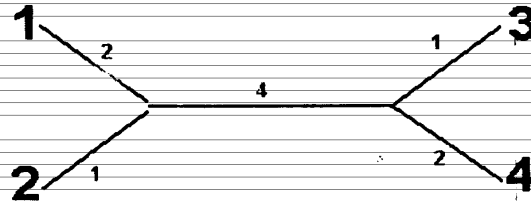
وبحساب البعد بين هذه التتابعات الأربع نحصل على النتائج التالية:

	1	2	3	4
1	-	3	7	8
2	-	-	6	7
3	-	-	-	3
4	-	-	-	-

وباستعمال هذه المعلومات يمكن تصميم شجرة قرابة نجمية (عديمة الجذر)

بين التتابعات الأربع حيث أن قيم الإستبدالات (الإختلافات) بين أزواج العينات

مقدرة كأطوال فروع الشجرة كما هو مبين كالتالي:



وعادة ما يفضل كثير من الباحثين التفكير في الشجرات التطورية على أن لها

سلف مشترك أو على أنها شجرة لها جذر rooted، لذلك لجأت كثير من الطرق لحسابات

خاصة تحقق ذلك.

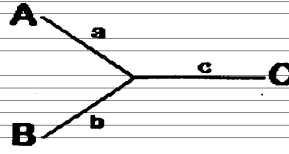
١.٢.٢.٧. طريقة Fitch & Margoliash.

وهي من أوائل الطرق التي اقترحت في هذا الصدد في أوائل السبعينات من

القرن الماضي والتي تحدد الوحدات داخل الشجرة في ثلاثيات threes ثم تحسب قيم

البعد لأذرع تلك الشجرة. وفيما يلي تعريف بأسس هذه الطريقة:

[١] تقترح شجرة من ثلاث وحدات كالتالي :



[٢] تحسب أطوال الأذرع جبرياً بناء على درجات البعد والمبينة كالتالي:

	A	B	C
A	-	22	39
B	-	-	41
C	-	-	-

Distance from A to B = $a + b = 22$ (1)

Distance from A to C = $a + c = 39$ (2)

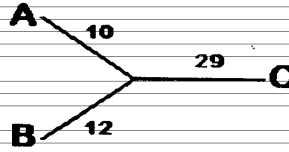
Distance from B to C = $b + c = 41$ (3)

للتخلص من المجهول تطرح (2) من (3) فيكون الناتج

$$(b + c) - (a + c) = 41 - 39 = b - a = 2 \quad \text{..... (4)}$$

بالجمع (1) و (4) الناتج سيكون $2b = 24$; $b = 12$ وبالتعويض في (1) و (2) تكون قيم

$10 = a$ و $29 = c$. وعلى ذلك تمثل الشجرة كالتالي:

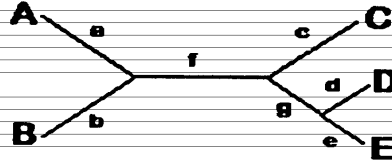


ولكن ما هو الحال إذا ما كان عدد العينات أكثر من ثلاث، مثلاً خمسة

A, B, C, D, E ، وعلاقات البعد بينهم يلخصها الجدول التالي:

	A	B	C	D	E
A	-	22	39	39	41
B	-	-	41	41	43
C	-	-	-	18	20
D	-	-	-	-	10
E	-	-	-	-	-

ويمكن تمثيل الشجرة كالتالي:



فيتم التعامل معها كالتالي:

[١] أولا يتم تحديد أكثر اثنين من العينات تحت الدراسة قرباً من بعضهما البعض من

جدول البعد الوراثي وفي هذه الحالة ستضح أنهما العينات D و E

[٢] نعيد حساب درجات البعد على أنها ثلاث عينات، فالحساب بعد D عن A, B, C

مجتمعة يأخذ متوسط البعد لهم الثلاث عن D، أي $(39 + 41 + 18)/3 = 32.7$

وللبعد بين E و A, B, C و يتابع نفس الطريقة سيكون 34.7 ، وعليه نعدل الجدول

ليصبح

	D	E	Ave. A,B,C
D	-	10	32.7
E	-	-	34.7
Ave. A,B,C	-	-	-

[٣] ويمكن أيضاً حساب العلاقة بين D و A,B,C وكذلك بين E و A,B,C عن طريق حساب

متوسطات أطوال الأفرع على الشجرة:

$$D \text{ to } E: d + e = 10 \quad (1)$$

$$D \text{ to } ABC: d + m = 32.7 \quad (2)$$

$$\text{where "m"} = g + (c + 2f + a + b)/3$$

$$E \text{ to } ABC: e + m = 34.7 \quad (3)$$

وبالطرح المعادلة (٢) من المعادلة (٢) فإن $d - e = -2$ وبالجمع مع المعادلة (١)

نحصل على أن $d = 4$; $2d = 8$; وبالتعويض فإن $e = 6$.

[٤] والآن نتعامل مع العينات D و E على أنها عينة واحدة ونبنى جدول جديد لقيم البعد

كالآتي:

	A	B	C	(DE)
A	-	22	39	40
B	-	-	41	42
C	-	-	-	19
(DE)	-	-	-	-

ومنها يمكن حساب قيم c ستساوي 9 وقيمة g ستساوي 5.

[٥] نعيد الخطوات السابقة حتى نحدد جميع قيم أطوال الأفرع على الشجرة تحت

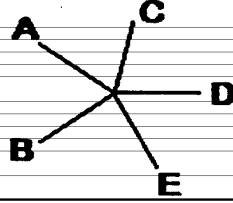
الدراسة.

ملخص خطوات طريقة Fitch-Margoliash:

- ١- حدد أكثر اثنين من العينات قربا.
- ٢- اعتبر باقي العينات وحدة واحدة، ثم احسب متوسط البعد بينهما ومتوسط العينات المجمعة.
- ٣- استعمل تلك العلاقات الجديدة في حساب أطوال أفرع الشجرة لهما.
- ٤- أعد الخطوات السابقة مع تجميع عينات أخرى ومنها احسب أطوال الأفرع وهكذا.

٢.٢.٢.٧ طريقة حساب ربط المتجاورات Neighbor-joining algorithm

أن هذه الطريقة قريبة الشبه جدا بالطريقة السابقة، فالعينات (التتابعات) التي يتم اختبارها للوصل بينها تعتمد على تقدير مربع الانحرافات least-square لأفضل التقديرات لطول أفرع الشجرة. وتبدأ الطريقة بتحديد الشجرة النجمية حيث لا توجد أى إتصالات بين العينات المتجاورة، كما هو موضح بالشجرة التالية:



وتحديد أى من العينات سيتم ربطها معا عن طريق حساب مجموع أطوال الأفرع

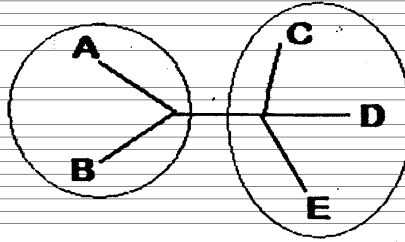
للشجرة بناء على المعادلة:

$$S_{mn} = \frac{\sum d_{im} + d_{in}}{2(N-2)} + \frac{d_{mn}}{2} + \frac{\sum d_{ij}}{N-2}$$

where ij represent all sequences except m and n , and $i < j$.

وعلى سبيل المثال دعنا نقول أن التتابعات (العينات) المتجاورة التي ستوصل أولا

هى A و B، كما هو مبين بالشجرة التالية:



٢.٢.٢.٧. طريقة UPGMA.

وهذه الطريقة هي أكثر الطرق استعمالاً لتقدير البعد و تحديد الشجرات التطورية مستخدمة تتابعات الـ DNA والبروتين. و هذا الأسـم UPGMA هو إختصار لطريقة الحساب المعروفة بأسم Un-weighted Pair Group Method with Arithmetic Mean. ويبدا الحساب بمحاولة تجميع clustering التتابعات في أزواج (الأقرب لبعضهما) وربطهما مع بعضهما فيما يسمى بالندبة أو النتوء node ثم زوج ثانى ثم ثالث وهكذا وهى ذلك نقوم ببناء الشجرة من أسفل إلى أعلى upwards فالنتوء الأول يعلوه الثانى والثانى يعلوه الثالث وهكذا وحواف (حدود) النتوء يحسب من فرق ارتفاع هذه الوحدات (النتوءات). والمسافة بين كل تجميعتين clusters مثلا C_i و C_j تحسب من المعادلة التالية،

$$d_{ij} = \frac{1}{|C_i| + |C_j|} \sum_{p \in C_i, q \in C_j} d_{pq}$$

where $|C_i|$ and $|C_j|$ are the number of sequences in clusters i and j , respectively

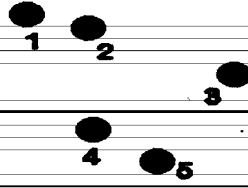
وحسابات البرنامج لطريقة UPGMA algorithm يمكن تلخيصها في الآتي:

- [١] أربط كل تتابع i لمجموعته C_i cluster.
- [٢] حدد كل ورقة (وحدة) من أوراق الشجرة T لكل تتابع، و وقعها عند الارتفاع صفر.
- [٣] حدد المجموعة من i و j بإستعمال المعادلة السابقة التى تعطى أقل تقدير لقيمة d_{ij} .
- [٤] حدد مجموعة ثانية k حيث $C_k = C_i \cup C_j$ ثم حدد d_{ki} مع كل التتابعات.
- [٥] حدد النتوء k بالنسبة للنتوء بن i و j ، و وقعها على إرتفاع يساوى $d_{ij} / 2$.
- [٦] أضف k للمجموعة وهكذا يستمر العمل لباقي التتابعات.

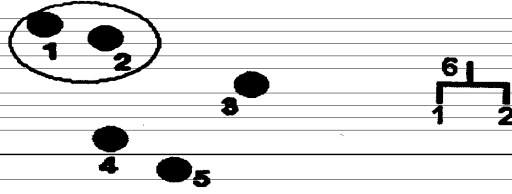
مثال تطبيقي على طريقة UPGMA :

إفترض وجود ٥ تتابعات يراد تحديد العلاقات بينها، وممثلة كنقاط على رسم

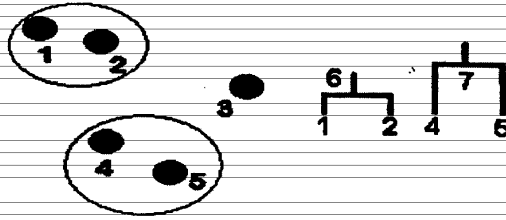
بيانى يحدد أماكنها، كالتالى :



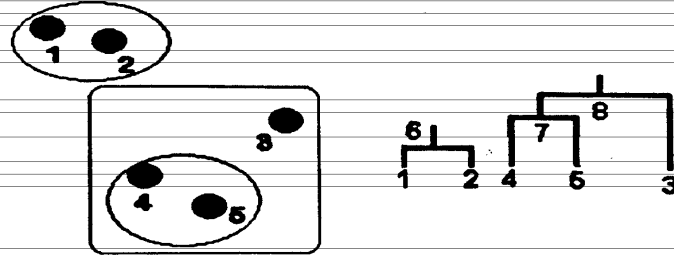
الآن إختار المجموعتين clusters الأكثر قربا من بعضهما البعض، وهما التتابعان ١ و ٢، كون مجموعة واحدة منهما (كما هو موضح بالشكل التالي) ثم أنشأ نتوء node يمثلها عند ارتفاع $d_{12}/2$:



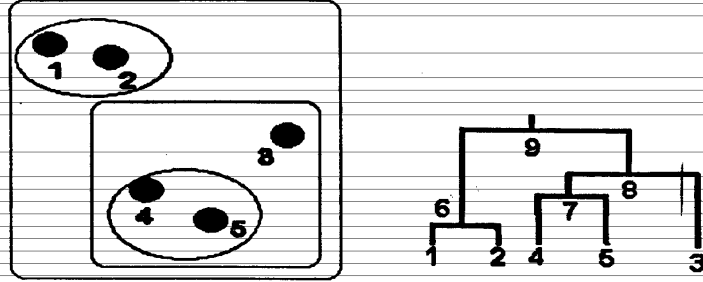
تابع العمل فإختار المجموعتين التاليتين من حيث القرب، وهما ٤ و ٥ وضمهما في مجموعة واحدة، وتابع تكوين الشجرة، كما هو مبين كالتالي:



والمجموعة التالية هي التي ستربط بين مجمعة (٥/٤) ومجموعة رقم ٣، كالتالي:



الآن لم يتبقى سوى مجموعتين لذا يجب الربط بينهما للإنتهاء من عمل الشجرة كالتالى:

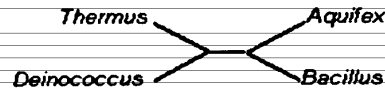


طرق تحديد الشجرات الفيلوجينية التى تناولناها بالشرح فى الأجزاء السابقة

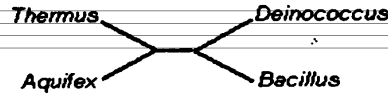
ومع أنها مستعملة على نطاق كبير فى التجارب الوراثية تعطى نتائج لأبأس بها إلا أن كفاءتها محدودة خصوصا فى الدراسات التطورية فعلى سبيل المثال هناك أربعة أقسام من

البكتيريا معروفة وراثيا وتطوريا باختلافها وهى:

- ١- قسم (أى وحدة تقسيمية) *Deinococcus*، وهى بكتيريا مقاومة للأشعة
- ٢- قسم *Thermus* وهى محبة للحرارة thermophilic، وهناك أدلة تجريبية قوية بأن هذان القسمان متقاربان وراثيا و ٣- القسم *Aquifex* وهى محبة
- أيضا للحرارة و ٤- القسم *Bacillus* وهى محبة للماء mesophilic، والقسمان الثالث والرابع لاتربطهما أى علاقة قرابة بالقسمين الأول والثانى، وعلى ذلك كان المتوقع أن تكون الشجرة بينهم كالتالى:



وباستعمال معظم الطرق السابقة ونتيجة التشابه في تتابعات الـ DNA بين *Aquifex* و *Thermus* فإنهما يقعان متجاورتان على الشجرة خلافاً للمتوقع كالآتي:



وللتغلب على تلك المشاكل يلجأ العلماء إلى طرق حسابية أخرى لتحديد الشجرة وتعرف طرق الأرجحية العظمى *maximum likelihood methods*، حيث أن هذه الطرق تعتبر أكثر كفاءة خصوصاً عند التعامل مع الأفرع الطويلة للشجرات والبيانات الغريبة مثل التتابعات الشاذة وعند إنحراف نسبة الإستبدالات *transitions/transversions* عن الواحد أو عند اختلاف معدل التحول في منطقة من التتابع عن مجاوراتها بمعنى أن يكون معدل التطور بطئاً في منطقة وسريع في منطقة أخرى. مع العلم أن طرق الأرجحية قد تؤدي لعمل شجرة خاطئة كالمثال السابق إلا أنها تختلف عن الطرق السابقة في قدرتها على عرض الشجرات الأقل أرجحية ولا ترفضها تماماً.

٢.٢.٧. طرق الأرجحية العظمى *Maximum Likelihood Methods*.

طرق الأرجحية العظمى هي طرق تعتمد على تقدير الاحتمال نسبة لنموذج *model* حسابي فرضي محدد. فعلى سبيل المثال إذا ما القينا بقطعة نقيود عادية (النموذج يحدد أن لها وجهان – الصورة والكتابة) فإن أرجحية الحصول على الصورة من رمية واحدة هي 0.5. وإذا نص النموذج أن قطعة النقيود هذه غير طبيعية وعلى وجهيها صورتين، فإن الأرجحية في هذه الحالة ستساوي الواحد الصحيح، أي أن طرق الأرجحية تعتمد على تقدير الاحتمال وفقاً لنموذج مفترض. وفي تصميم النموذج *model* تبعاً

لطرق الأرجحية لتعيين شجرة فيلوجينية لتتابعات عينات ما فإن النموذج يجب أن يشتمل على جزئين وهما:

- ١- المكون composition، ويختص بنسب النيوكليوتيدات الأربع (A,C,G,T) لبعضها البعض
- و ٢- النهج process ، ويختص بمعدلات الاستبدالات substitutions بين هذه النيوكليوتيدات الأربع. وفيما يلي تعريف مبسط بالأساس الحسابي لبرامج تقدير الشجرات بطرق الأرجحية العظمى.

فلحساب أرجحية تتابع ما (هل نيوكليوتيدة واحدة من A للتبسيط) فهذا لن توجد شجرة وسنستبعد أيضا جزء النهج process من النموذج وسيتوقف فقط على جزء المكون component، فإذا كان النموذج ينص على أن التتابع ١٠٠٪ A إذن الأرجحية ستكون واحد، وإذا ما نص النموذج بأنه ١٠٠٪ C إذن الأرجحية ستكون صفر وإذا ما نص النموذج بأن A تساوى ٢٢٪ فإن الأرجحية ستساوى ٠,٢٢.

وما هي الأرجحية لو كان التتابع من عدة نيوكليوتيدات، هنا لابد من تضمين النموذج جزء النهج process حيث أن الإستبدالات بين القواعد ستلعب دورا، والذي ستمثله المصفوفة التالية:

$$P = \begin{bmatrix} 0.976 & 0.01 & 0.007 & 0.007 \\ 0.002 & 0.983 & 0.005 & 0.01 \\ 0.003 & 0.01 & 0.979 & 0.007 \\ 0.002 & 0.013 & 0.005 & 0.979 \end{bmatrix}$$

حيث أن ترتيب القواعد في المصفوفة سيكون أبجديا أى A ثم C ثم G ثم T، ومن تلك القيم فإن احتمال بقاء A دون تغير هو 0.976 وأن تستبدل بـ C هو 0.01 أو تستبدل بـ G هو 0.007 وهكذا لباقى الإحتمالات الستة عشر. وجزء النموذج الخاص بالمكون composition والتي سنرمز لها بـ π ، دعنا نقول أن نسبهم ستكون $\pi = [0.1, 0.4, 0.2, 0.3]$

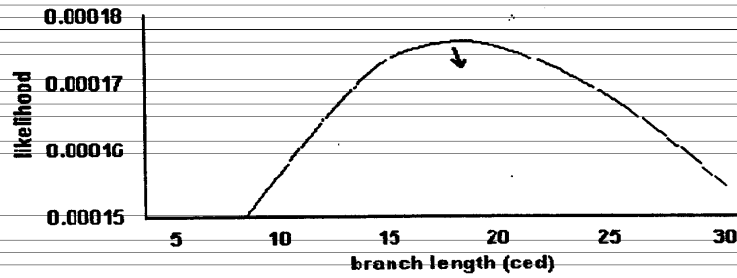
وبفرض وجود تتابعان هما ss و CCGT فسوف يربطهما فرع واحد بسيط (الشجرة)، إذن أرجحية هذا الفرع يمكن أن تحسب من العلاقة التالية:

$$= \pi_c P_{c \leftarrow c} \pi_c P_{c \leftarrow c} \pi_g P_{g \leftarrow g} \pi_t P_{t \leftarrow t} = 0.4 \times 0.983 \times 0.4 \times 0.983 \\ \times 0.1 \times 0.007 \times 0.3 \times 0.79 = 0.0000300$$

وحيث أن فرع الشجرة السابق بسيط، فيجب حساب فروع الشجرة عندما تكون متغيرة الطول (كما هو الحال في معظم الأحيان) وذلك نسبة لوحدات زمنية تطورية (للسهولة يفترض أنها وحدات بعد تطوري إعتبارية certain evolutionary distance [ced]، وهي في مثالنا السابق قدرها وحدة واحدة. ولكن ما التغيرات المنتظرة لمصفوفات النهج إذا ما كانت قيم ced أكبر من واحد (تضرب المصفوفة في نفسها بعدد قيم وحدات الزمن التطوري)، كما في المصفوفات التالية:

$$P^2 = \begin{bmatrix} 0.953 & 0.02 & 0.013 & 0.015 \\ 0.005 & 0.966 & 0.01 & 0.02 \\ 0.007 & 0.02 & 0.959 & 0.015 \\ 0.005 & 0.026 & 0.01 & 0.959 \end{bmatrix} \& P^3 = \begin{bmatrix} 0.93 & 0.029 & 0.019 & 0.022 \\ 0.007 & 0.949 & 0.015 & 0.029 \\ 0.01 & 0.029 & 0.939 & 0.022 \\ 0.007 & 0.038 & 0.015 & 0.94 \end{bmatrix}$$

ويمكننا ملاحظة أن زيادة طول الفرع يؤدي بالضرورة إلى خفض احتمالات أن تبقى القواعد دون تغير (الصف القطري on diagonal) بينما احتمالات الإستبدالات تزيد (الصفوف غير القطرية off diagonal). والرسم البياني المبين بشكل (٤-٧) يوضح العلاقة بين وحدات البعد التطوري وقيم الأرجحية.



شكل (٤-٧) : العلاقة البيانية بين وحدات البعد التطوري وقيم الأرجحية المحسوبة.

ويتضح أن قيم الأرجحية تصل إلى حدها الأقصى عند وحدات بعد تطوري تتراوح بين ١٥ — ٢٠ ced ، كذلك إذا ما زادت قيم البعد بقدر كبير فإن مصفوفة النهج ستؤول إلى قيم المكون π ، كما هو موضح بالمصفوفة التالية:

$$P^{10^6} = \begin{bmatrix} 0.1 & 0.4 & 0.2 & 0.3 \\ 0.1 & 0.4 & 0.2 & 0.3 \\ 0.1 & 0.4 & 0.2 & 0.3 \\ 0.1 & 0.4 & 0.2 & 0.3 \end{bmatrix}$$

وعادة تستعمل المصفوفات السابقة في صورة لوغاريتمية logarithmic لصغر الأرقام ولحاولة فصل المنهج عن المكون في النموذج، كما أنها تتيح لنا حساب طول الفرع نسبة إلى معدل الاستبدال لكل موقع per site substitution خلافاً للوحدات الإعتبارية السابقة (ced) كذلك يتيح لنا التعامل مع أطوال للفروع من الصفر حتى ما لانهاية، لذلك تصبح المصفوفة كالآتي:

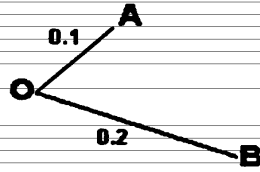
$$\log P = \begin{bmatrix} -0.0244 & 0.0101 & 0.0067 & 0.0076 \\ 0.0025 & -0.0176 & 0.005 & 0.0101 \\ 0.0034 & 0.0101 & -0.021 & 0.0076 \\ 0.0025 & 0.0134 & 0.005 & -0.021 \end{bmatrix}$$

حيث يلاحظ أن مجموع الصفوف سيؤول للصفر في هذه الحالة. ويجب تحويل المصفوفة للأساس الأسى حتى تعبر عن احتمالات منسوبة للإستبدالات لكل موقع كالآتي:

$$scaled\ P = \begin{bmatrix} -0.122 & 0.05 & 0.034 & 0.038 \\ 0.05 & -0.353 & 0.101 & 0.202 \\ 0.034 & 0.101 & -0.21 & 0.076 \\ 0.038 & 0.202 & 0.076 & -0.315 \end{bmatrix}$$

ولحساب العلاقة بين تتابعين مختلفين (A & B) عند أطوال أفرع مختلفة، كما

هو مبين كالتالي:



حيث تتابع A هو CCAT وتتابع B هو CCGI وتم ربطهما بمنشأ غير معروف (O)، وباستعمال المصفوفات لعدلات إستبدال قدرها $P_{0.1}$, $P_{0.2}$, $P_{0.3}$ يمكن حساب الأرجحية بعدة طرق، أبسطها، تعتمد على أن المسافة الكلية بينهما هي 0.3 (من A إلى O ومن O إلى B) إذن تستعمل احتمالات المصفوفة $P_{0.3}$ ، ومنها نحصل على التقدير التالي:

$$\pi_c P_{c-c} \pi_c P_{c-c} \pi_a P_{a-a} \pi_i P_{i-i} = 0.4 \times 0.786 \times 0.4 \times 0.786 \times .1 \times 0.08 \times 0.3 \times 0.747 = 0.000177$$

٢.٧. برامج الكمبيوتر للعلاقات الفيلوجينية.

Software for phylogenetic relations.

حاولنا في الأجزاء السابقة التعريف بالأسس الحسابية والإحصائية لطرق تحديد العلاقات الفيلوجينية لعينات من تتابعات النيوكلووتيدات أو البروتينات لتفهم الأسس التي بنيت عليها تلك الطرق مع بيان مميزات وعيوب كل منها. ومما سبق يتضح أن هذه الطرق تعتمد على أساليب رياضية معقدة وصعبة وإن تنفيذها يدويا يعتبر متعب للغاية ويحتاج لزمان طويل، لذلك فإن تلك الحسابات تتم اعتمادا على برامج عديدة للكمبيوتر لتسهيل المهمة وتوفير الوقت. وخلال السنوات القليلة الماضية توفر عدد كبير من تلك البرامج المتخصصة وصل عددها الآن لأكثر من ٢٠٠ برنامج مختلف تتباين فيما بينها من حيث طبيعة اللغات والأنظمة المستخدمة، ومنها ما هو متاح للجميع بدون مقابل (خصوصا الإصدارات القديمة منها) وأغلبها يباع لمن يرغب. وعادة يمكن تبويب تلك البرمجيات تبعا للفرض الذي صممت من أجله، ويمكن تلخيص بعض منها كالتالي:

PHYLP -1

وهو من أشهر البرامج في الدراسات الفيلوجينية، ويمكن الحصول عليه

مجانياً من الإنترنت <http://evolution.gs.washington.edu/phylip.html>

<http://evolution.gs.washington.edu/phylip.html>

وهو متوافق مع معظم الأنظمة العالمية للكمبيوتر، وعن طريقه يمكن القيام

بحسابات الطرق المقتصة parsimony ومصفوفات البعد distance matrix والأرجحية

العظمى maximum likelihood وغيرها مستعملاً بيانات من تتابعات الـ DNA أو الـ RNA

أو البروتينات أو مواقع القصر أو أي بيانات غير مستمرة (1/0) أو التكرارات الجينية. وهو

(مع البرنامج التالي) يمثل أكثر من ٨٠٪ من بيانات القرابة الفيلوجينية المنشورة في

الدوريات العالمية المتخصصة.

PAUP -2

وهي حزمة من البرامج وتعنى Phylogenetic Analysis Using Parsimony حيث

صمم في الأصل لهذا الغرض ثم طور فيما بعد ليقوم بكافة الحسابات التي يقوم بها

البرنامج السابق، وهو ليس بالمجان وإن كان ثمنه غير باهظ (حوالي ٩٠ دولار لنظام

Windows). ولزيد من التفاصيل عن البرنامج يمكن اللجوء إلى الموقع:

<http://paup.csit.fsu.edu/>

MacClade -3

يعتبر من البرامج الرائدة في تحليل العلاقات التطورية مستعملاً عدد متباين

من طرز البيانات منها الصفات المظهرية غير المستمرة والبيانات الجزيئية، ويباع بحوالي

١٢٥ دولاراً، ولزيد من التفاصيل يمكن الاطلاع على الموقع:

<http://phylogeny.arizona.edu/macclade/macclade.html>

Hennig86 -4

وهو برنامج سريع لتقدير المعايير المقتصة وطرق للبحث عن أنسب التفرعات

للشجرات الفيلوجينية، والبرنامج مشفر وللحصول على الشفرة لابد من دفع اشتراك

ومصاريف البريد. والبرنامج يمكنه التعامل مع حوالي ٨٠ عينة taxa وحوالي ٩٩٩ صفة

مختلفة، ولزيد من التفاصيل عنه يمكن الرجوع إلى:

Farris, J.S. 1989, Hennig86: a PC-DOS program for phylogenetic analysis. *Cladistics* 5: 163.

Random Cladistics -5

وهو برنامج للشجرات يعتمد في تقديراته على برنامج Hennig86 وهو قادر عن البحث عن الجزر بين الشجرات، والوصف المفصل للبرنامج موجود على الموقع:
<http://research.amnh.org/~siddall/rc.html>

حيث يمكن الحصول على نسخة منه بلغة الـ DOS.

AutoDecay -6

وهو برنامج خاص يمكنه أن يبني معاملات التداعيات decay indices اعتمادا على الشجرات المبنية على أحد البرامج الأخرى وهو برنامج PAUP. وهو مكتوب بلغة Perl ومتوافق مع أغلب أنظمة الكمبيوتر المتاحة، ويمكن الحصول عليه من شبكة المعلومات من الموقع:

http://www.bergianska.se/index_forskning_soft.html.

DNA Stacks -7

وهو برنامج لتحليل تتابعات الـ DNA ورصها ولا يقوم ببناء العلاقات الفيلوجينية، وهو مساعد للبرامج الأخرى مثل PAUP و PHYLIP ولكنه صالح لنظام Macintosh فقط. ويمكن الحصول عليه من شبكة المعلومات من الموقع:

<http://biology.fullerton.edu/deernisse/dnastacks.html>.

TreeRot -8

وهو برنامج لعمل الشجرات اعتمادا على الإحصائيات المقتصدة ولكن متوافق مع نظام Macintosh ولزيد من التفاصيل يمكن الرجوع إلى الموقع:

<http://people.bu.edu/msoren/TreeRot.html>

RA -9

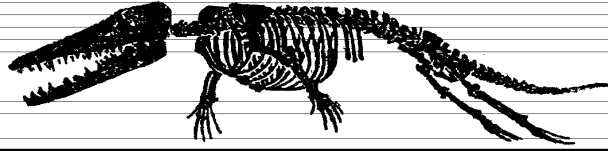
وهو برنامج صغير لتحليل بيانات من تتابعات النيوكليوتيدات وهو متوافق مع العديد من الأنظمة ويعتبر رخيص الثمن (حوالي ٣٠ دولار).

٤.٧. تطبيقات الفيلوجينيا Applications of phylogenetics.

الفيلوجينيا هي محاولات لدراسة التاريخ التطوري لمجموعة من الكائنات تمثل وحدات تقسيمية taxa ما على أسس وراثية. وحديثا تركزت هذه الدراسات على فحص التتابعات الجينية أو البروتينية (الجزئيات الحيوية) لذلك قد تعرف أيضا باسم الجينومكس المقارن comparative genomics. والفيلوجينيا تختص بدراسة وتفهم موضوعين رئيسيين ، أولهما بناء الشجرات التطورية مع تحديد تفرعاتها المختلفة ، وثانيهما هو إستغلال تلك الشجرات لمعرفة مسارات تطورية جديدة، كذلك لا يجب أن نغفل الهدف التقليدي للفيلوجينيا ألا وهو التقسيم classification & systematics. وفيما يلي بعض من أهم تطبيقات تلك الدراسات.

١.٤.٧. التطور Evolution.

تطور الكائنات موضوع واسع وكبير، لكن من الأمثلة الشيقة في هذا المجال هي الدراسات المستفيضة التي أجريت لمعرفة أسلاف الحيتان whale ancestors . فقد ظل موضوع تطور الحيتان والدرافيل (cetaceans) من العضلات غير المفهومة في علم التطور لزم من طويل، فكيف يمكن تفسير أن هذا الحيوان الثديي mammal الكبير قد إرتد إلى البحر مرة أخرى وتأقلم على العيش في تلك البيئة المائية على عكس مسار التطور المعروف بأن الكائنات البرية terrestrial قد خرجت من البحر منذ ملايين السنين خلال مسارها التطوري ! لكن حديثا وفي أواخر القرن الماضي (١٩٨٠ – ١٩٩٥) تم إكتشاف عدد من الحفريات التي تمثل حلقات الوصل بين الثدييات الحافرية artiodactylans والحيتان البحرية، فقد اكتشف Gingerich أول تلك الحفريات في باكستان وسميت *Packcetus* ثم حفريات الحوت ذو الأرجل *Ambulocetus* في باكستان ومصر بواسطة Thewissen، وشكل (٥-٧) يبين هذا الحيوان المنقرض.



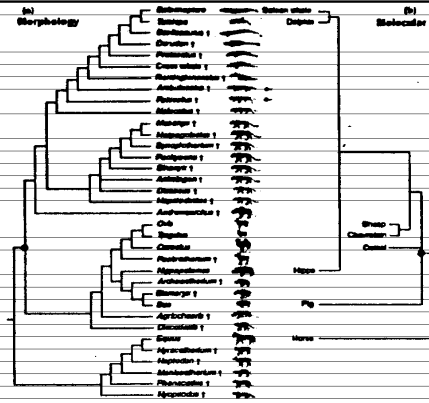
شكل (٥-٧) ، حفرية الثدي المنقرض من حوال ٤٠ مليون عام المعروف باسم *Ambulocetus*.

وببناء الشجرة الفيلوجينية للثدييات الحافرية اعتماداً على الصفات المظهرية لما هو باق منها وللمجموعة الحفريات fossils لحيوانات تابعة لها ولكنها انقرضت، يتضح أن الحيتان والدرافيل قد نشأت من سلف مشترك (ثديى حافرى) يمثل حلقة الوصل بين الحيوانات البرية والمائية يعرف باسم *Ambulocetus*. ولكن عندما بنيت الشجرة باستعمال البيانات الجزيئية من تتابعات الـ DNA المستخلص من الميتوكوندريا والأنوية mt DNA & nuclear DNA لعدد من ذات الحافر، إتضح أن الحوت أكثر قرابة إلى فرس النهر (سيد قشطة) ويمثلان فصيلتان شقيقتان sister clade، وأن فرس النهر أقرب للحوت أكثر من قرابته لباقى ذات الحافر مثل الأبقار و الحصان وغيرها، كما هو موضح بشكل (٦-٧).

٧.٤.٢. أدلة الطب الشرعي Forensic medicine evidence.

أصبح استعمال تقنيات البيولوجيا الجزيئية وخصوصاً تكنيك البصمة الوراثية DNA fingerprinting فى مجالات الطب الشرعي فى قضايا البنية وإثبات النسب وكذلك للتعرف على المتهمين فى الجنايات من الأمور شائعة الاستعمال منذ حوالي ربع قرن مضى من الزمان. ولكن ما يهمنا فى هنا هو إمكانية استعمال الشجرات الفيلوجينية المعتمدة على تحليل تتابعات النيوكليوتيدات فى مجالات الطب الشرعي أو الطب الوقائى.

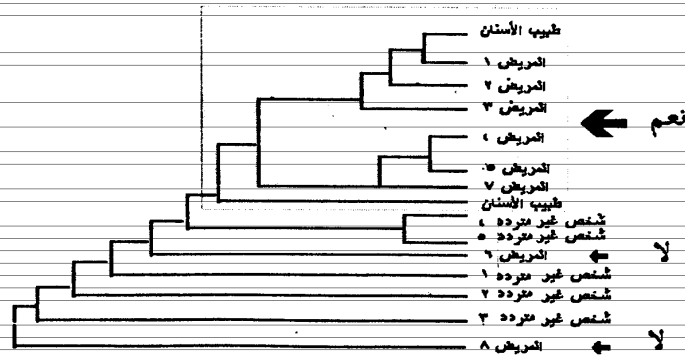
ففى دراسة شيقة قام بها كل من OU وزملاؤه سنة ١٩٩٢ لتحديد مصدر الإصابة بفيروس مرض الإيدز بين مجموعة من المرضى المترددين على عيادة أحد أطباء الأسنان بالولايات المتحدة، قاما بعزل تلك الفيروسات من هؤلاء المرضى وطبيب الأسنان ومن عدد آخر من الأفراد المصابين بالإيدز ولكن غير مترددين على عيادة هذا الطبيب (عينة ضابطة)، وقاما



شكل (٦-٧) : شجرات فيلوجينية للتشخيصات ذات العاقر (a) بناء على الصفات المظهرية morphology والحفريات (b)

بناء على تناهات الـ DNA الميتوكوندري والنووي (molecular). حيوانات منقرضة = †.

ببناء شجرة فيلوجينية اعتماداً على تحليل تناهات النيوكليوتيدات لتلك العزلات من الفيروسات، كما هو موضح بشكل (٧-٧). فإمكانهم تحديد المرضى الذين أصيبوا بالمرض عن طريق عيادة هذا الطبيب والمرضى الذين أصيبوا من مصادر أخرى.




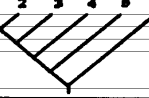
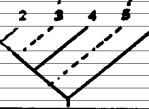

شكل (٧-٧) : شجرة فيلوجينية لتناهات فصائل فيروس الأيدز.

٧.٤.٧. التنبؤ بوظائف الجينات Gene function prediction.

خلال السنوات القليلة الماضية ظهرت الحاجة الملحة للتعرف على وظائف الجينات وتحديدها وذلك مواكبتاً لازدهار دراسات الجينومكس في العديد من الكائنات، وخلال تلك الدراسات تم عزل تنابعات (جينات) لم تكن معروفة من قبل وكان على الوراثيون أن يتعرفوا أو على الأقل أن يحاولوا التنبؤ بوظائفها. وفي سبيل ذلك إتبع العلماء طرق عديدة ولكن أكثرها شيوعاً هو محاولة دراسة درجات التشابه similarity بين تلك الجينات مجهولة الوظيفة وشبيهاتها (بناءً على صف alignment تنابعات النيوكليوتيدات – الباب السادس) من الجينات معلومة الوظيفة ثم محاولة التنبؤ بوظيفة هذا الجين مجهول الوظيفة. وهذه الطرق أثبتت نجاحاً لا بأس به في التعرف على وظائف العديد من الجينات، ولكن في بعض الحالات فإن التشابه الجزئي لم يعكس بالضرورة التشابه الوظيفي، يمكن الرجوع إلى Hillis سنة ١٩٩٤ لتحديد بعض من تلك الحالات النادرة. لذلك يفضل العديد من الوراثيين إستعمال أكثر من طريقة لهذا الغرض وعدم الاكتفاء بدراسة التشابه الجزئي. خصوصاً أن مفهوم التماثل الوراثي homology غير محدد، فعلى سبيل المثال الجدول التالي يلخص عدد من مصطلحات التماثل التي قد تتداخل في مفاهيمها مع بعضها البعض:

المصطلح	المعنى
Homolog	الجينات المتشابهة التي يربطها سلف واحد مشترك (مثل كل جينات الجلوبين globins في الثدييات).
Ortholog	الجينات المتشابهة التي انفصلت عن بعضها diverged بعد حادثة (خطوة) تطورية ما speciation event (مثل جيني β -globin في الإنسان والفأر).
Paralog	الجينات المتشابهة التي انفصلت عن بعضها diverged بعد حدوث تضاعف جيني gene duplication (مثل جيني β -globin و γ -globin في الإنسان).
Xenolog	الجينات المتشابهة التي انفصلت عن بعضها البعض diverged بعد حدوث إنتقال جزئي لأحد الجينات lateral gene transfer (مثل جينات مقاومة المضادات الحيوية في البكتيريا).

ويعتقد العلماء أن طرق الصف لتحديد التماثل غير قادرة على عكس البعد التطوري لهذه الجينات لذلك اقترح Eisen سنة ١٩٩٨ طريقة جديدة تجمع بين المعلومات الفيلوجينية (الشجرات) وطرق الصف التقليدية للتنبؤ بوظائف الجينات والتي أطلق عليها اسم "الفيلوجينوميكس" Phylogenomics. وشكل (٧-٨) يلخص أهم خطوات هذه الطريقة.

	<p>1 حدد الجين غير مطوم للوظيفة</p> <p>2 تعرف على مجموعة الجينات المذكورة</p>
	<p>3 بناء على صف التتابعات alignment حدد الشجرة الفيلوجينية لهذه المجموعة من الجينات</p>
	<p>4 وقع على الشجرة وظائف الجينات المعروفة ما أمكن حيث الجينات ١ و ٢ حدد لها وظيفة واحدة والجينات ٤ و ٦ حدد لها وظيفة واحدة لغرض بناء الجين ٣ فغير مطوم الوظيفة مثله مثل الجين تحت الدراسة رقم ٥</p>
	<p>5 يمكن التنبؤ بأن وظيفة الجين ٥ قد تكون مشابهة للجينات ٤ و ٦ بينما الجين ٣ فوظيفته أقرب للجينات ١ و ٢</p>

ومع افضلية هذه الطريقة ونجاحها في التعرف على وظائف العديد من الجينات مجهولة الوظيفة إلا أنها غير مؤكدة بالقدر الكافي فمازال هناك حالات لا يمكن التعرف عليها لذلك وجب البحث على طرق أكثر دقة وتحديداً.

•

•

•

•

•

٨. البروتيومكس

Proteomics

إعداد : أحمد الشهاوى

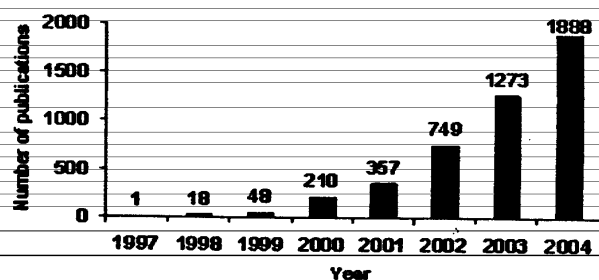
٨.١. مقدمة.

بعد الإنتهاء من مشاريع قراءة التتابع الكامل للجينوم العديد من الكائنات الحية مثل جينوم نبات الأرابيدوبسيس *Arabidopsis* ، جينوم الإنسان Human genome ، الخ، بدأ إهتمام العلماء بعصر مابعد الجينوم post-genomic era. هذا يعنى كيفية الإستفادة من المعلومات الجينومية الهائلة التى تم تجميعها وتخزينها لفهم الوظائف والأنظمة الخلوية والإختلافات بينها. ففى الوقت الحالى توجد فجوة كبيرة فى الوسائل والطرق لربط المعلومات الهائلة المتاحة عن التراكيب الجينومية المختلفة ووظائف البروتينات المشفرة من هذه المعلومات الوراثية (ربط تركيب الجينوم بوظيفته).

بالنظر إلى بعض الجينومات التى تم الإنتهاء من تحديد تتابعها يمكن توقع التحديات الضخمة التى تنتظر العلماء فى عصر مابعد الجينوم. فعلى سبيل المثال، فإن جينوم خميرة الخباز (*Saccharomyces cerevisiae*) والذى تم الإنتهاء من قراءته، فإن حوالى ٢٠% من الـ Open Reading Frames (ORFs) مازال غير معروف وظيفتها وتحتاج إلى طرق تعتمد على تحليل ودراسة الجينوم ككل. إن تركيب وتنظيم الجينومات الأكثر تعقيدا مثل جينوم الإنسان تجعل المهمة أكثر صعوبة مقارنة بجينوم الخميرة. وعند الأخذ فى الإعتبار تنظيم التعبير الجينى على مستوى النسخ، وما بعد النسخ، والترجمة، وما بعد الترجمة وكذلك مرونة الجينوم وقابليته للتغير فى حدود معينة، فإن المشكلة تتضخم للدرجة لا يمكن تصورها ويصبح تصور ودراسة وفهم وظيفة الجينوم مهما كان صغيرا فى حجمه مهمة شبه مستحيلة. وتصبح كل المحاولات لدراسة وظيفة الجينوم هى محاولات لفهم جزئى لوظيفة الجينوم. لذلك يحتاج عصر مابعد الجينوم لطرق تحليل الجينوم بالكامل whole system-based analysis لمحاولة فهم لوظيفته الجينوم ككل. ومن أهم الوسائل المتاحة لعصر مابعد الجينوم هو التحليل الكمى لتعبير

بعض الجينات على مستوى النسخ quantitative monitoring of expression level لعدد كبير من الجينات معا من خلال استخدام microarray technology. كذلك يمكن دراسته وظيفته حين معين من خلال دراسة التغيرات في التعبير الجيني differential gene expression، ولكن في حالة الكائنات المعقدة فإن المعلومات التي يتم الحصول عليها من دراسات التعبير الجيني على مستوى النسخ لا تعكس التغيرات التي تحدث على مستوى البروتين في معظم الأحيان. هذا يرجع إلى طرق التنظيم المختلفة والعديد التي تحدث أثناء الترجمة وما بعد الترجمة translational and posttranslational regulations. فقد أشارت دراسات عديدة على ضعف العلاقة بين تعبير الـ RNA والبروتين، وبالتالي فمن الصعوبة استخدام مستوى النسخ (RNA) لتوقع مستوى الترجمة (البروتين)، وتحورات ما بعد الترجمة post-translational processing. ففي دراسة على جينوم الخميرة، دلت النتائج على أنه لا بد من وجود طرق أدق لدراسة التعبير الجيني للجينوم ككل. وقد شملت هذه الدراسة التعليم بواسطة العوامل الوراثية المتنقلة transposon tagging لعدد ٣٦ جين من جينات الإنقسام الميوزي meiosis لدراسة تعبيرهم الجيني على مستوى البروتين، لكن لم يتمكنوا إلا من تحديد وظيفة ١٧ جين فقط من هذه الجينات عند دراستها بتقنية الـ microarray. من ذلك يمكن القول بأن العلاقة بين التعبير الجيني على مستوى الـ RNA، والتعبير الجيني على مستوى البروتين ليست فقط علاقته معقدة ولكنها قد تكون مضللة.

لا يوجد في الوقت الحالي طرق عامه لتحليل وفهم الوظائف البيولوجية للجينات. هذه الطرق يجب أن تتضمن العديد من طرق تحليل الجزيئات البيولوجية المختلفة المشتركة في الوظائف الخلوية في نفس الوقت. تعتبر البروتينات من أهم الجزيئات البيولوجية التي تقوم بالوظائف الخلوية بطريقه مباشره. ولقد تم استخدام مصطلح البروتيومكس proteomics بواسطة Wilkins et al. سنة ١٩٩٦ لوصف دراسة كل البروتينات المنتجة بواسطة الجينوم ككل وفي صورته كميته عند وقت معين أو استجابته لظروف معينه. هذا المجال مازال في مراحله الأولى ويحتاج إلى جهد كبير لتطويره في المستقبل، وبالرغم من أن هذا المجال مازال في مراحل تطوره الأولى فإنه يمثل نقطة جذب للباحثين كما هو موضح في الشكل التالي لعدد الأبحاث المنشوره والتي تحتوى على كلمه بروتينومكس (شكل ٨-١).



شكل (١-٨)، عدد الأبحاث المنشورة في الفترة من ١٩٩٧ حتى ٢٠٠٤. تم الحصول على هذه الأعداد بإجراء بحث في الـ PubMed باستخدام كلمة proteomics.

٢.٨. عزل وفصل البروتينات.

بالرغم من أن مجال البروتيومكس proteomics حديث نسبياً فإن بعض تقنياته مستخدمة من حوالى ٢٥ سنة. مثل الفصل الكهربى للبروتينات فى اتجاهين 2D gel electrophoresis الذى يستخدم لفصل ودراسة البروتينات منذ أكثر من ٢٥ سنة. كذلك فإن استخدام Mass Spectroscopy مع الفصل الكهربى للبروتينات فى اتجاهين 2D gel electrophoresis كونا أول جيل من نظم دراسة البروتيومكس first generation of proteomic platforms. و أثناء السنوات القليلة الماضية تم تطوير العديد من الطرق لفصل وتنقية البروتينات والتي تعتبر الأساس للجيل الأول من نظم دراسة البروتيومكس وفيما يلى ملخص لهذه الطرق.

١.٢.٨. عزل البروتين Protein Isolation.

أول وأهم خطوة فى دراسة البروتيوم proteome (المعين البروتينى) هى إستخلاص البروتين من البيئة الموجود بها. ففى هذه الخطوة يجب الإهتمام بعزل البروتيوم proteome فى أقرب صورة كان موجود عليها فى الخلية وتقليل تأثير طريقه الفصل. يجب التأكيد على أنه لا توجد طريقه عامه لفصل البروتينات نظراً لاختلاف

طبيعة العينات المستخدمة ولكن توجد طرق مختلفة لفصل البروتينات بطريقة مرضية من العينات المختلفة. لفصل بروتينات السيتوسول cytosolic proteins يتم تحليل الخلايا بطريقة بسيطة واستخدام الـ supernatant الناتج. هناك العديد من الطرق المختلفة لعزل البروتينات وهي متاحة على الموقع التالي بشبكة المعلومات الدولية (<http://expasy.cbr.nrc.ca/ch2d/protocols>)، وبعض منها قد تكون بسيطة كاستخدام صدمة الضغط الأسموزي osmotic shock، أو طرق أعنف مثل استخدام الموجات فوق صوتية ultrasonication. إن اختيار الطريقة المستخدمة لفصل البروتين تتوقف على سهولتها وملاءمتها لفصل البروتينات في اتجاهين (تركيز الملح والمذيبات الأيونية). وبالرغم من أن معظم هذه الطرق مدروسة بطريقة جيدة إلا أن فصل بعض مكونات البروتيوم proteome مثل البروتينات الغير محبة للماء hydrophobic proteins مازال من الصعوبة بمكان. ففي السنوات الأخيرة تم تطوير بعض الكيماويات والطرق لتحسين فصل البروتينات الغير محبة للماء hydrophobic proteins، ومع هذا التقدم النسبي في طرق فصل البروتينات للتحليل ثنائي الاتجاه إلا أن نتائج دراسات البروتيوم تشير إلى أن الاختلافات الخلوية في العينة المستخدمة وكذلك الاختلاف في أماكن تخزين البروتين داخل الخلايا يؤثر في نتائج هذه الدراسات.

٨.١.٢.١. العينة المتجانسة Homogenous Sample.

الخلايا المنزرعة cell cultures ربما تكون من أفضل مصادر البروتين للدراسة البروتيوم حيث أنها تمثل مجموعة من الخلايا المتجانسة وتوفر بروتينوم متحكم فيه إلى حد ما. كذلك فإن ظروف نمو وتحلل الخلايا يمكن التحكم فيه والذي يقلل من الاختلافات بين المزارع والاختلافات غير الحقيقية artifacts.

٨.١.٢.٢. العينة الغير متجانسة Heterogeneous Sample.

العينات غير المتجانسة مثل عينات الأنسجة تجعل تحليل البروتيوم لها عملية معقدة للغاية وذلك بسبب الاختلافات الكبيرة في طريقة الحصول على العينات وطريقة إستخلاص البروتين. فمثلاً، عينات الأنسجة التي تستخدم في دراسات البروتيوم في الإنسان يتم تجميعها من مرضى من المستشفيات. حيث لا توجد عملية لرقابه كافيه على جمع وحفظ هذه العينات والذي يمثل مصدر كبير للاختلافات غير الحقيقية في

البروتيوم نتيجة الإستجابة للمؤثرات الخارجيه. عينات الأنسجه غير متجانسه بطبيعتها (تحتوى على أنواع خلايا مختلفه بأحجام مختلفه وبنسب خلايا مختلفه) وهذا يؤدي إلى وجود إختلافات فى البروتيوم من نفس النسيج أو بين الأنسجه المختلفه. لذلك فإن النتائج المستخلصه من مثل هذه الدراسات يجب الوصول اليها بعد دراسة عدد كبير من العينات. فى الواقع فإن البروتيوم المستخلص من عينات الأنسجه يتكون من عديد من البروتيوم كل منها ناتج من مجموعه من الخلايا فى العينه أو النسيج. وفى هذه الحاله فإن البروتيوم المختلط الناتج لايمكن مقارنته حتى بدراسة عديد من العينات.

٨.١.٢.٢ مجاميع الخلايا داخل العينه Sub-Population of Cells.

أصبح من الضرورى فصل عشره من الخلايا من العينه البيولوجيه لتحليلها وتحديد المعلومات المستمده من كل بروتيوم. حديثا تم إستخدام أنواع متقدمه من اجهزه فصل الخلايا cell sorters وطرق مناعيه متقدمه لفصل تحت عشائر الخلايا لتحليل البروتيوم.

٨.١.٢.٤ التراكيب تحت الخلويه Subcellular components.

يمكن أيضا تبسيط البروتيوم إلى مجموعه البروتينات المستخلصه من أحد مكونات الخليه حيث أن نقل وتوزيع البروتينات إلى عضيات معينه داخل الخليه عمليه معروفه فى مجال بيولوجيا الخليه cell biology. فيمكن فصل أحد عضيات الخليه والتعامل مع بروتينات هذه العضيه على أنها بروتيوم مستقل.

٨.٢.٢ الفصل الكهربى للبروتينات فى الإتجاهين.

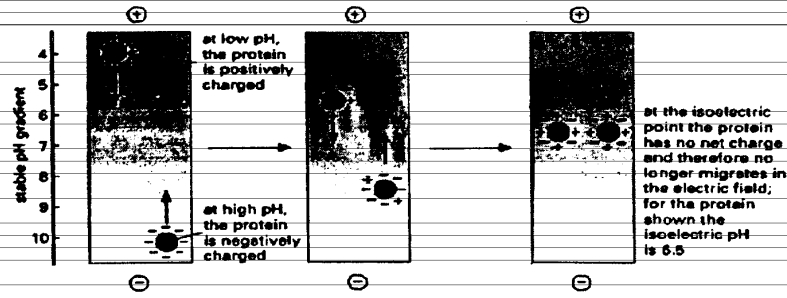
Separation by 2D electrophoresis.

بعد عزل وتنقيه البروتينات يلزم فصلها بطريقه معينه لفصل مكونات مخلوط البروتين إلى بروتينات فرديه. فى مجال البروتيومكس، فإن طريقه الفصل المفضله هى الفصل الكهربى فى إتجاهين 2D electrophoresis. هذه الطريقه يمكن أن تستخدم لفصل عده الألف من البروتينات فى تجربه واحده.

٨.٢.١. أساسيات فصل البروتينات في إتجاهين.

Principles of 2D separation.

بدأ الإستخدام الأمثل للفصل الكهربى للبروتينات في إتجاهين بعد تطور شرائط التدرج في الـ pH الثابتة (IPG) Immobilized pH Gradient في الثمانينات من القرن الماضى، والذي يتم بيلمره مجموعه من الجزيئات الأحادية monomers والتي تحمل وظيفه الـ ampholytes في وجود acrylamide. لذلك بتغيير تركيز هذه الجزيئات على طول الشريط ينتج تدرج ثابت في الـ pH. وتوجد تجاريا اشراطه تحتوى على تدرج pH في مدى مختلف لفصل البروتينات في إتجاهين. الإتجاه الأول للفصل (تركيز البروتينات تجاه نقطه التعادل الكهربى) Isoelectric Focusing. . للفصل في هذا الإتجاه، فإن مستخلص البروتين يستخدم لتشبيح وإعادة تمدد شريط الـ IPG قبل الفصل. هذا يضمن أن البروتين سيكون موزع بشكل متساوى على الشريط ويمنع ترسيب البروتين لزياده تركيزه في مناطق معينه على الشريط. يتم بعد ذلك إمرار تيار كهربى تدريجى خلال الشريط. البروتينات المشحونه بشحنه موجب (أى انها تقع في منطقه من الشريط ذات pH أقل من درجه التعادل الكهربى pI لها) فإنها تتحرك ناحية القطب السالب cathode وتواجه بدرجات pH أعلى حتى تصل إلى درجات التعادل الكهربى pI لها وتصبح عندها متعادله وتتوقف عن الحركة. البروتينات المشحونه بشحنه سالبه (أى انها تقع في منطقه من الشريط ذات pH أعلى من درجه التعادل الكهربى pI لها) فإنها تتحرك ناحية القطب الموجب anode وتواجه بدرجات pH أقل حتى تصل إلى درجات التعادل الكهربى pI لها وتصبح عندها متعادله وتتوقف عن الحركة. النتيجة النهائيه للفصل، فإن البروتينات التى لها نفس نقطه التعادل الكهربى فإنها تركز وتوجه باستمرار إلى منطقه من الشريط ذات درجه pH مماثله لدرجة تعادلها الكهربى pI (شكل ٨-٢).



شكل (٨-٢) : فصل البروتين بالتركيز تجاه نقطة التعادل الكهربى. عند درجات الـ pH المنخفضة (تركيز مرتفع لأيونات الهيدروجين)، فإن البروتين يحمل شحنة موجبة.

الإتجاه الثانى للفصل، هو فصل البروتينات على أساس وزنها الجزيئى SDS-PAGE، حيث يتم صب وبلمره محلول الأكريلاميد بين قطعتين من الزجاج يفصل بينهما فاصل يتراوح سمكه من ١.٥-١ مم. يتم وضع الشريط الناتج من الفصل فى الإتجاه الأول أفقيا على SDS-PAGE (SDS-Polyacrylamide Gel Electrophoresis). يتم مرور تيار كهربى ويتم فصل البروتينات بناء على وزنها الجزيئى. أنظمة فصل البروتينات فى الإتجاه الأول والثانى متوفرة تجاريا من مصادر مختلفة.

٢.٢.٢.٨. تحديد وإظهار البروتينات المفصولة فى إتجاهين.

Resolution of Separated Proteins.

توجد طرق مختلفة لتحديد وإظهار البروتينات المفصولة فى إتجاهين. هذه الأنظمة تعتمد على تعليم البروتينات قبل أو بعد الفصل أو ربما بعد الفصل فى الإتجاه الأول. وأهم هذه الطرق هى: (١) تعليم البروتينات قبل الفصل Pre-Labeling of proteins ويندرج تحتها التعليم بالعناصر المشعة Radiolabeling، حيث يتم تعليم البروتينات قبل الفصل فى الإتجاه الأول عادة داخل الخلية *In vivo* باستخدام العناصر المشعة radioisotopes. هذا التعليم يتم باستخدام جزئ معلم (حامض أمينى) مثل ^{35}S -methionine فى مزارع الخلايا، حيث يتم إظهار البروتينات المعلمة بعد ذلك بتعريض الجيل لفيلم X-ray أو شاشة phosphor imager كما فى

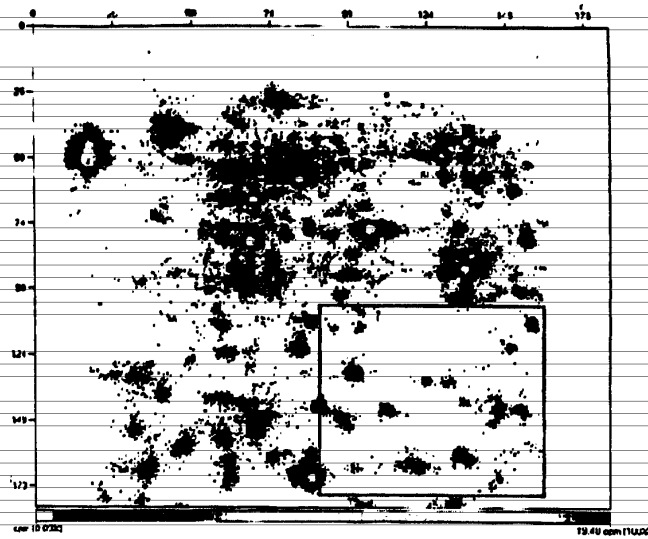
شكل (٨-٣). هذا النوع من التعليم غير مناسب لدراسة عينات الحيوان والإنسان. ويندرج تحتها أيضا طرق التعليم الفلورسنتى Fluorescent labeling حيث يتم التعليم بصبغات فلورسنتيه متعادله كهربيا والتي ترتبط بالبروتين قبل الفصل الكهربى فى الإتجاه لأول. هذه الصبغات يعاب عليها أن حساسيتها أقل ولكن فى الفترة الأخيرة تم تطوير صبغات يصل حد إيضاها مستوى النانوجرام من البروتينات. إن كل هذه الصبغات الفلورسنتيه تحتاج إلى نظام خاص لتحديد شدة الفلورسنس fluorescence لكان البروتين على الجيل.

توجد طرق أخرى تسمى multiphoton detection (MPD) وتعتمد على تعليم البروتين *In vivo* بواسطة كميات قليلة من اليود المشع ¹²⁵I أو ¹³¹I. هذا النظام يمكنه إظهار كميه من البروتين على مستوى الأتومول (amol) (٨-٤) فى هذه الحالة يتم إستخدام كميه من البروتين فى مستوى النانوجرام. شكل (٨-٤) يوضح الفصل فى إتجاهين حوالى ١٥٠ نانوجرام من مستخلص بروتين بكتريا *E. coli* والتي تم تعليمها قبل الفصل بواسطة اليود المشع ¹²⁵I. أما مجموعة طرق التعليم الثانية فتتميز بكونها طرق لتعليم البروتينات بعد الفصل Post-Labeling of Proteins، فتعليم البروتينات بعد الفصل هى الطريقة الشائع لإظهار البروتينات بعد الفصل الكهربى فى إتجاهين ويتم ذلك بواسطة صبغات مثل blue Coomassie (حدود الإظهار بهذه الطريقة Limit of Detection (LOD) هى ٥٠ نانوجرام من البروتين) أو الصبغ بالفضه Silver Staining (حدود الصبغ بهذه الطريقة حوالى ٥ نانوجرام). توجد طرق أخرى لإظهار البروتينات معتمده على الصبغ العكسى reverse staining. بالرغم من أن الصبغ بالفضه Silver Staining يحتاج إلى عمل أكثر إلا أنه يعتبر من أكفأ التكنيكات وأكثرها حساسيه لإظهار البروتينات على الجيل (فى الجيل). شكل (٨-٥) يوضح جيل gel مصبوغ بالفضه مقصول فى إتجاهين لحوالى ١٥٠ ميكروجرام بروتين. ويلاحظ أن كل منطقه مصبوغه spot تحتوى على بروتين أو أكثر ويمكن إستخدامها لتحليل ودراسة هذه العينات.

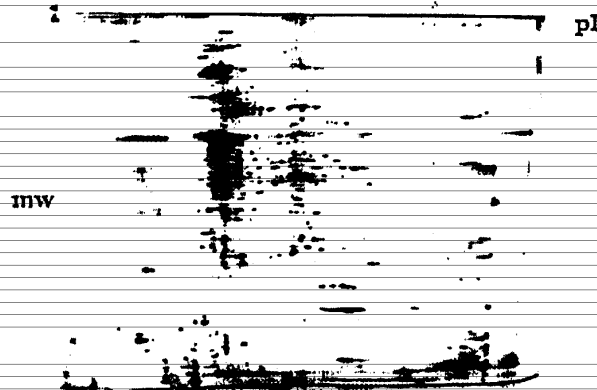


شكل (٨-٢) : الفصل الكهربى للبروتين الكلى لـ *E. coli* فى اتجاهين معلمه بمخلوط من الأحماض الأمينية المشعة قبل الفصل. الإتجاه الأفقى يوضح الفصل على حسب نقطة التعادل الكهربى *pI* والإتجاه الراسى يمثل الفصل على حسب الوزن الجزيئى.

كذلك توجد تكتيكات أخرى أكثر حساسية تعتمد على التعليم الفلورسنتى. هذه الطرق لا تظهر البروتينات فى صورة مرئية على الجيل *gel* ولكنها تعتمد على أجهزة تستخلص البروتينات من الجيل وتوصفها مباشرة. هذا الفصل ثنائى الإتجاه يعتبر من أكفأ الطرق لفصل البروتينات والذى يمكننا من فصل أكثر من ١٠٠٠٠ بروتين على الجيل الواحد.



شكل (٤-٨)، الفصل الكهربائي لبروتينات بكتريا *E. coli* (١٠٠ ميكروجرام) في اتجاهين مصبوغه قبل الفصل باليود المشع ¹²⁵I ، الاتجاه الأفقي يوضح الفصل على حسب نقطة التبادل الكهربائي *pI* والاتجاه الرأسي يمثل الفصل على حسب الوزن الجزيئي.



شكل (٤-٩)، الفصل الكهربائي لبروتينات *Sulfolobus solfataricus* (١٥٠ ميكروجرام) في اتجاهين مصبوغه بعد الفصل بالفضة *Silver staining*. الاتجاه الأفقي يوضح الفصل على حسب نقطة التبادل الكهربائي *pI* والاتجاه الرأسي يمثل الفصل على حسب الوزن الجزيئي.

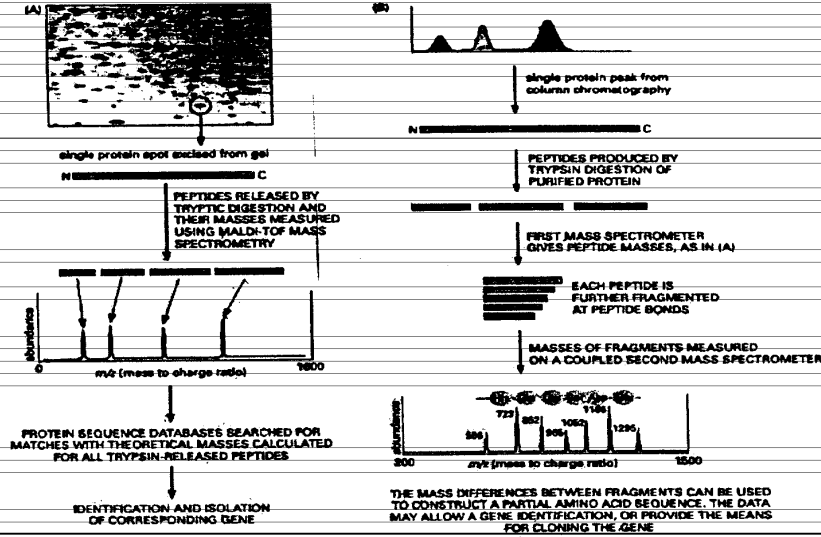
٨.٢. التطورات الحديثة في دراسة البروتيومكس. New Developments.

تشير الدراسات السابقة عن البروتيومكس أن مستوى إظهار وتحديد البروتينات المفصولة هو مستوى الفمومول fmol. وهذا يعتبر مستوى مرتفع بالنسبة للبروتينات ذات التركيز المنخفض. وتعتبر هذه المشاكل الصعبة بالنسبة للبروتينات منخفضة التعبير والتي قد تتسبب في انخفاض النشاط الإنزيمي enzyme kinetics أثناء الهضم في وجود تركيز قليل من مادة التفاعل (البروتين). كذلك تؤدي إلى انخفاض حساسية الـ mass spectrophotometer تحت هذه الظروف. وحديثاً بالذکر إن صعوبة تحديد وإظهار البروتينات ذات التركيز المنخفض على الجيل ثنائي الإتجاه 2D gel يمكن إرجاعها إلى وجود بروتينات أخرى ذات تركيز عالٍ تتداخل مع إظهار تلك البروتينات shadowing of low abundant proteins by high abundant proteins. تبعاً لذلك تم تطوير طرق أخرى حديثة للتغلب على بعض هذه المشاكل. البروتينات عالية التركيز تحجب البروتينات ذات التركيز الأقل على الجيل حيث أن الأول ينتج مناطق صبغ أكبر large spots. ويبدو أن زيادة حجم الجيل قد تقدم حلاً لهذه المشكلة. بالإضافة إلى استخدام شرائط IPG تحتوي على مدى ضيق من الـ pI (zoom 2D gel) تم تصميمها للتغلب على حجب البروتينات ذات التركيز العالي للبروتينات ذات التركيز الأقل protein shadowing. في هذا التكنيك يتم استخدام مجموعة شرائط IPG تحتوي على مدى قصير ومختلف من الـ pI مع كل عينة بروتين. فمثلاً، بدلاً من استخدام شريط واحد IPG لتغطية الـ pI من ٢-١٠ طوله ٨ سم فإنه يتم استخدام من ٤-٥ شرائط كل منها ٨ سم كل منها يغطي حوالاً ١,٥ من مدى الـ pI. ينتج من ذلك عديد من الـ 2D gel الكبيرة كل منها يغطي مدى معين من الـ pI والتي يقل فيها تأثير البروتينات ذات التركيز العالي على إظهار البروتينات قليلة التركيز.

٨.٢.١. تحليل الكتلة الضوئي. Mass Spectrometry.

التطورات والتقدم الحديث في الـ mass spectrometry مكنت من الإظهار السريع وتحديد تتابع الأحماض الأمينية الجزئي لكميات قليلة من البروتينات. هذه الطرق تعتمد على هضم البروتينات المفصولة من الـ 2D gel إما في الجيل أو بعد

نقل هذه البروتينات إلى أغشيه. الهضم في الجيل in-gel digestion وهو الشائع وذلك لسهولة. هذه الطريقة تعتمد على قطع البروتين من الجيل يدويا أو أوتوماتيكيا وهضمها بالترسين ثم تحديد نتائجها من الأحماض الأمينية جزئيا. ويتم استخدام نموذج الببتيدات الناتجة من الهضم أو التتابع الجزئي من الأحماض الأمينية لتحديد الجين الذي ينتج هذا البروتين شكل (٨-٦). التطورات الحديثة في مجال البروتيومكس تشمل تطبيقات عديدة في مجالات مختلفة من أهمها مجال دراسة الأساس الجزيئي للأمراض واكتشاف أدوية جديدة لعلاج تلك الحالات.



شكل (٨-٦): شكل تخطيطي يوضح استخدام طريقة تحليل الكتلة للتعرف على البروتينات وتحديد

تتابع الأحماض الأمينية لمكوناتها الببتيدية.

تستخدم هذه الطريقة لتحديد كتلة البروتينات والببتيدات الناتجة بعد تحليلها الإنزيمي بدقه واستخدام هذه المعلومات للبحث في قواعد البيانات عن الجينات الخاصة بهذه البروتينات. حيث يتم فصل البروتين من الجيل ثنائى الإتجاه وهضمه بالترسين ثم

قياس كتله الببتيدات الناتجه من الهضم بواسطه جهاز قياس الكتله . يتم بعد ذلك البحث فى قواعد البيانات عن جينات تنتج بروتينات لها نفس نموذج الهضم بالتريسين.

كذلك يستخدم تحليل الكتله أيضا لتحديد تتابع الأحماض الأمينية للببتيدات

الناتجه من الهضم بالتريسين. فى هذه الحاله فإن البروتينات المفصولة بالـ

chromatography كبروتينات مفردة يتم هضمها بالتريسين ثم يتم تحديد كتله

الببتيدات الناتجه بجهاز تحديد الكتله بنفس الطريقه كما فى شكل (٨-١٦). لتحديد

التتابع الدقيق للأحماض الأمينية فإن كل ببتيده ناتجه يتم تحديد تتابعها من الأحماض

الأمينيه بتكسير روابطها الببتيديه. هذا ينتج مجموعه من الببتيدات المختلفه فى الطول

بحامض أمينى واحد. يتم نقل هذا الخليوط من الببتيدات إلى جهاز تحليل كتله آخر

متصل لتحديد كتلتها. الإختلاف فى الكتله بين الببتيدات المختلفه فى حامض أمينى

واحد تستخدم لتحديد الحامض الأمينى المفقود. بترار هذه العمليه يمكن تحديد تتابع

الأحماض الأمينية لجزء من البروتين الكلى (شكل ٨-٦٦).

1

2

3

4

5

6

٩. المراجع References

- أبويوسف، أميرة يوسف و أحمد يوسف المتينى – التقنية الحيوية الجزيئية وتحسين النباتات – منشأة المعارف – الإسكندرية – ٢٠٠٢.
- المتينى، أحمد يوسف – تنظيم الفعل الجينى فى الكائنات الراقية – منشأة المعارف – الإسكندرية – ١٩٩٧.
- المتينى، أحمد يوسف – مدخل الوراثة الجزيئية – منشأة المعارف – الإسكندرية – ١٩٩٤.

1. Adams, M. D., M. B. Soares, A. R. Kerlavage, C. Fields, and J. C. Venter, 1993. Rapid cDNA sequencing (expressed sequence tags) from a directionally cloned human infant brain cDNA library. *Nature Genetics*, 4: 373 – 380.
2. Antonopoulos, J. E. *Genomics*. Xlibris Corporation, 2000.
3. Barnes, M., and I. C. Gray (Eds.). *Bioinformatics for Geneticists*. John Wiley & Sons, 2003.
4. Benson, D. A., I. Karsch-Mizrachi, D. J. Lipman, J. Ostell, and D.L. Wheeler (2005). GenBank. *Nucleic Acids Res.*, 33: 231-244.
5. Bergeron, B. P. *Bioinformatics Computing*. Prentice Hall, 2002.
6. Bishop, M. J., and C. J. Rawlings (Eds.). *DNA and Protein Sequence Analysis: A Practical Approach* (Practical Approach Series, No 171), IRL Press. 1996.

7. Bishop, M. J. (Ed.). *Genetic Databases*. Academic Press, San Diego, 1999.
8. Bornholdt, S., and H. G. Schuster (Eds.). *Handbook of Graphs and Networks : From the Genome to the Internet*. Vch. Verlagsgesellschaft MbH., 2003.
9. Brown, S. M. *Bioinformatics: A Biologist's Guide to Biocomputing and the Internet*. Eaton Pub.Co., 2000.
10. Campbell, A. C., and L. J. Heyer. *Discovering Genomics, Proteomics, and Bioinformatics*. Benjamin Cummings Pub., 2003.
11. Causton, H. C., J. Quackenbush, and A. Brazma. *Microarray Gene Expression Data Analysis: A Beginner's Guid.*, Blackwell Publishers, 2003.
12. Collado-Vides, J., and R. Hofstadtt (Eds.). *Gene Regulation and Metabolism: Post-Genomic Computational Approaches* (Computational Molecular Biology), MIT Press, 2002.
13. Davidson, Erick H. *Genomic Regulatory Systems: Development and Evolution*. Academic Press, San Diego, 2001.
14. Durrett, R., and R. Durrett. *Probability Models for DNA Sequence Evolution*. Springer Verlag, 2002.
15. Eisen, J. A., 1998, Phylogenomics: improving functional prediction for uncharacterized genes by evolutionary analysis. *Genome. Res.*, 8: 163-167.
16. Fall, C., E. Marland, J. Wagner, and J. Tyson (Eds.). *Computational Cell Biology*. Springer Verlag, 2002.

17. Felsenstein, J. *Inferring Phylogenies*. Sinauer Associates, 2003.
18. Fitch, W.M., and Margoliash, E., 1968, The construction of phylogenetic trees. II. How well do they reflect past history? *Brookhaven Symp. Biol.*, 21(1):217-242.
19. Fitch, W.M., and Margoliash, E., 1967, Construction of phylogenetic trees. *Science*, 155(760):279-284.
20. Fogel, B., and D. W. Corne (Eds.) *Evolutionary Computation in Bioinformatics*, Morgan Kaufmann. 2002.
21. Gingerich, P.D., Haq, Mu., Zalmout, I.S., Khan, I.H., and Malkani, M.S., 2001, Origin of whales from early artiodactyls: hands and feet of Eocene Protocetidae from Pakistan. *Science*, 293(5538): 2239-2242.
22. Griffiths, A. J..F.; Gelbart, W. M.; Miller, J. H.; Lewontin, R. C. *Modern Genetic Analysis*: W. H. Freeman & Co., 1999.
23. Gribskov, M., and J. Devereox (Eds.). *Sequence Analysis Primer*. Oxford University Press, 1992.
24. Grotewold, E. (Edt). *Plant Functional Genomics* (Methods in Molecular Biology, Vol. 236) (Methods in Molecular Biology (Clifton, N.J.), V. 236.), 2003.
25. Hall, B. G. *Phylogenetics Trees Made Easy: A How-To Manual for Molecular Biologists*. Sinauer Associates, 2001.
26. Hartemink, A. J., D. K. Giford, T. S. Jaakkola, and R. A. Young. Using graphical models and genomic expression data to statistically validate models of genetic

regulatory networks," in Paciuc Symposium on Biocomputing, vol. 6, 2001.

27. Hawley, R. S., and M. Y. Walker. *Advanced Genetic Analysis: Finding Meaning in the Genome*. Blackwell Publishers, 2003.
28. Hedrick, P. W. *Genetics of Populations*. Jones & Bartlett Pub., 2000.
29. Heiskanen, M., Hellsten, E., Kallioniemi, O.P., Makela, T.P., Alitalo, K., Peltonen, L., and Palotie, A., 1995. Visual mapping by fiber-FISH. *Genomics*. 1;30(1):31-36.
30. Hillis, D. M. In *Homology: the hierarchical bases of comparative biology* (B. K. Hall, Ed.), pp. 339-368, Academic Press, San Diego 1994.
31. Hunt, S., and F. Livesey (Eds.). *Functional Genomics: A Practical Approach* (The Practical Approach Series, 235), Oxford Univ. Press, 2000.
32. James, P. (Ed.). *Proteome Research: Mass Spectrometry (Principles and Practice)*, Springer Verlag, 2001.
33. Kimura, M., and N. Takahata (Eds.). *Population Genetics, Molecular Evolution, and the Neutral Theory*. Selected Papers. University of Chicago Press. 1964.
34. Koonin, E. V., and Galperin, M. Y. *Sequence - Evolution - Function: Computational Approaches in Comparative Genomics*. Kluwer Academic Publishers, 2002.
35. Koski, T., and T. Koskinen. *Hidden Markov Models for Bioinformatics*, Kluwer Academic Publishers. 2001.

36. Klug, W., and M. R. Cummings. *Concepts of Genetics*, Fifth Edition, Printic Hall Inc., 1999.
37. Lesk, A. M. *Introduction to Bioinformatics*. Oxford University Press. 2002.
38. Letovsky, S. (Ed.). *Bioinformatics: Databases and Systems*. Kluwer Academic Publishers. 1999.
39. Levine, M. and Eric H. Davidson, 2005. Gene regulatory networks for development. Proc. Natl. Acad. Sci., USA, 102: 4936-4942.
40. Liebler, D. C. (Ed.). *Introduction to Proteomics: Tools for the New Biology*. Humana Press, 2001.
41. Nei, M. ,1996, Phylogenetic analysis in molecular evolutionary genetics. Annu. Rev. Genet., 30:371-403.
42. Nei, M., and S. Kumar. *Molecular Evolution and Phylogenetics*. Oxford Univ Press, 2000.
43. Ou, C.Y., Ciesielski, C.A., Myers, G., Bandea, C.I., Luo, C.C., Korber, B.T., Mullins, J.I., Schochetman, G., Berkelman, R.L., Economou, A.N.,1992, Molecular epidemiology of HIV transmission in a dental practice. Science, 256(5060):1165-1171.
44. Palzkill, T. *Proteomics*. Kluwer Academic Publishers, 2002.
45. Percus, J. K. (Ed.). *Mathematics of Genome Analysis*. Cambridge Univ. Press. 2002.
46. Rashidi, H. (Ed.). *Bioinformatics Basics: Applications in Biological Science and Medicine*. CRC Press. 1999.
47. Rastrigin, L. *This Chancy, Chancy, Chancy World*. Mir Publishers, 1973.

48. Ross, S. *A First Course in Probability* (6th Edition), Prentice Hall, 2002.
49. Rothstein, M. A. (Ed.). *Pharmacogenomics: Social, Ethical, and Clinical Dimensions*. Wiley-Liss, 2003.
50. Saccone, C., and G. Pesole. *Handbook of Comparative Genomics: Principles and Methodology*. Wiley-Liss, 2003.
51. Salemi, M., and Anne-Mieke Vandamme (Eds.). *The Phylogenetic Handbook: A Practical Approach to DNA and Protein Phylogeny*. Cambridge University Press, 2003.
52. Sankoff, D., and J. H. Nadeau. *Comparative Genomics - Empirical and Analytical Approaches to Gene Order Dynamics, Map Alignment and the Evolution of Gene Families*. Kluwer Academic Pub., 2000.
53. Sensen, C. W. (Ed.). *Essentials of Genomics and Bioinformatics*. John Wiley & Sons. 2002.
54. Sorensen, D. and D. Gianola. *Likelihood, Bayesian and MCMC Methods in Quantitative Genetics*. Springer Verlag, 2002.
55. Suhat, S. (Ed.). *Genomics and Proteomics: Functional and Computational Aspects*. Plenum Pub. Corp., 2000.
56. Till, B., Steven H. Reynolds, Clifford Weil, Nathan Springer, Chris Burtner, Kim Young, Elisabeth Bowers, E., C. A. Codomo, L. C. Enns, A. R. Odden, E. A. Greene, L. Comai and S. Henikoff , 2004. Discovery of induced point mutations in maize genes by TILLING. *BMC Plant Biology*, 4:12-17.

57. Thewissen, J.G., Williams, E.M., Roe, L.J., and Hussain, S.T., 2001, Skeletons of terrestrial cetaceans and the relationship of whales to artiodactyls. *Nature*, 413(6853):277-281.
58. Venter, J. Craig *et al.*, 2001, The Sequence of the Human Genome. *Science*, 1291 (5507):1304-1351.
59. Wang, J. T. L., Shapiro, B. A., Shasha, D. E. (Eds.). *Pattern Discovery in Biomolecular Data: Tools, Techniques, and Applications*. Oxford Univ Press. 1999.
60. Waterman, M. S. (Ed.). *Mathematical Methods for DNA Sequences*. CRC Press, 1989.
61. Wren, B., and N. Dorrell (Eds.). *Functional Microbial Genomics*. (Volume 33). Academic Press, 2003.

1

2

3

4

5

6

7

8

9

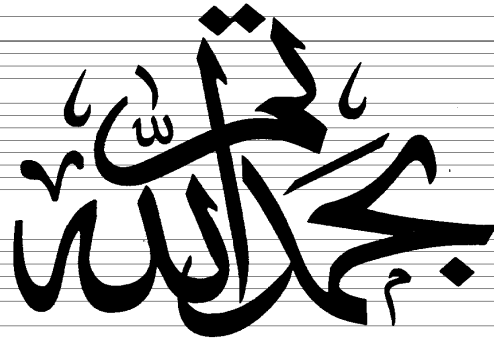
10

11

12

13

14



مكتبة بلستار المعرفة
لطباعة ونشر وتوزيع الكتب
كفر الدوار - الحدائق - بجوار نقابة التطبيقيين
٠١٢١١٥١٢٢٧ & ٠٤٥/٢٢٢٤٢٢٨

8
9
10
11

12
13
14
15

16
17
18